



مرکز تحقیقات اسلامی

اصفهان

گامی



الحق
علیه
صلاوة
وسلام

www.

www.

www.

www.

Ghaemiyeh

.com

.org

.net

.ir



مدیریت منابع اطلاعاتی وب

جلد اول

مبانی و تجربه‌های جهانی



به کوشش

دکتر اطلاعاتی ممتاز

فرزانه شادان پور

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

مدیریت منابع اطلاعاتی وب

نویسنده:

غلامعلی منتظر

ناشر چاپی:

سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران

ناشر دیجیتال:

مرکز تحقیقات رایانه‌ای قائمیه اصفهان

فهرست

۵	فهرست
۱۴	مدیریت منابع اطلاعاتی وب جلد ۱
۱۴	مشخصات کتاب
۱۵	اشاره
۲۰	فهرست مطالب
۲۲	سخن نخست
۲۴	به جای مقدمه
۲۶	فصل اول: مبانی مدیریت و آرشیو وب
۲۶	اشاره
۲۷	چکیده
۲۸	بایگانی شبکه وب: مباحث و روش ها
۲۸	اشاره
۲۸	۱-مقدمه
۲۹	۲-حفاظت از میراث
۳۱	۱-۲-به اندازه کافی خوب نیست؟
۳۷	۲-۲-۱-خود بقای رسانه؟
۴۰	۲-۳-۱-یک وظیفه غیر ممکن
۴۶	۳-ویژگی های وب برای حفاظت
۴۷	۱-۳-کاردینالیته وب
۵۰	۲-۳-وب به عنوان سیستم نشر فعال
۵۳	۳-۳-وب به عنوان یک محصول فرهنگی
۵۶	۴-روش های جدید برای رسانه جدید
۵۷	۱-۴-حفاظت از وب و زیر ساخت های اطلاعات
۵۹	۲-۴-فرآهم آوری
۶۱	۱-۲-۴-بایگانی جانی سرویس گیرنده
۶۷	۲-۲-۴-بایگانی تراکشی
۷۰	۲-۲-۴-بایگانی سرور- جانی
۷۳	۴-۳-سازماندهی و ذخیره سازی
۷۳	۱-۳-۴-نظام فایل های محلی به خدمت گرفته شده بایگانی ها
۷۳	توصیف
۷۵	توضیح
۷۸	استفاده ترجیحی
۷۸	ابزارها
۷۸	۲-۳-۴-بایگانی های مبتنی بر خدمت وب
۷۸	اشاره
۸۰	توصیف
۸۰	توضیح
۸۳	استفاده ترجیحی
۸۵	ابزارها
۸۵	۲-۳-۴-بایگانی غیروب

۸۵	توصیف
۸۷	توضیح
۸۷	مرجع
۸۷	اشاره
۸۷	۴-۳-۴ خلاصه
۸۸	۴-۴ کیفیت و تمامیت (کامل بودن)
۹۴	۵-بررسی عمومی مراحل اولیه جاری
۹۵	۱-۵-بازیگران پایگانی
۹۷	۲-۵-حوزه (دامنه)
۹۷	۵-۲-۱ پایگانی مرکزی سایت
۹۷	۵-۲-۲ پایگانی مرکزی عنوان
۹۹	۳-۲-۵ پایگانی مرکزی حوزه ای
۱۰۱	۳-۵-روش های استفاده شده
۱۰۳	۶-نتیجه گیری
۱۰۴	منابع
۱۲۳	چکیده
۱۲۴	از آرشیو اینترنت تا آرشیو در اینترنت
۱۲۴	اشاره
۱۲۴	مقدمه
۱۲۵	پیشینه: جاب اینترنتی اولیه
۱۲۵	آغاز به کار آرشیو اینترنت
۱۲۹	ساختار پیوندی و روپات های نواری
۱۳۱	۱۹۹۸:حضور داده های آرشیوی بر روی (تقریباً) هر دستکتاب
۱۳۱	۱۹۹۹: از نوار تا دیسک، یک خزشگر جدید و تصاویر متحرک
۱۳۳	۲۰۰۰: ایجاد مجموعه های موضوعی وب
۱۳۶	۲۰۰۱: دسترسی از طریق ماشین Wayback. آرشیو یازده سپتامبر
۱۴۰	۲۰۰۲: کتابخانه اسکندریه، کتابخانه سیار، و حق مؤلف
۱۴۳	۲۰۰۳: گسترش دستیابی ما به کتابخانه های ملی و مؤسسات آموزشی
۱۴۵	۲۰۰۴: آرشیو اروپا و پتاپاکس
۱۴۵	اشاره
۱۴۵	آینده
۱۴۵	منابع
۱۴۷	چکیده
۱۴۸	کاربرد وب و مطالعات مربوط به آن
۱۴۸	اشاره
۱۴۸	خلاصه
۱۴۸	اشاره
۱۴۹	۱-تحلیل محتوا
۱۵۱	۲-بررسی ها
۱۵۲	۳-تحلیل بلاغی
۱۵۴	۴-تحلیل گفتنمان

۱۵۵	۵-تحلیل دیداری
۱۵۷	۶-قوم نگاری
۱۵۸	۷-تحلیل شبکه
۱۵۹	۸-ملاحظات اخلاقی
۱۶۰	۹-نتیجه گیری
۱۶۰	منابع
۱۶۶	چکیده
۱۶۷	خصوصیات وب ایران: مریم پیروزمند
۱۶۷	اشاره
۱۶۷	۱.مقدمه
۱۷۰	۲.کارهای مرتبط در داخل و خارج
۱۷۲	۳.سامانه خودکار ارزیابی وب ایران
۱۷۸	۴.نتایج بدست آمده
۱۷۹	۴.۱.آمار وب گاه ها
۱۸۴	۴.۲.آمار صفحات
۱۸۸	۴.۳.آمار نوع صفحات
۱۹۱	۵.نتیجه گیری و کارهای آینده
۱۹۳	منابع
۱۹۳	اشاره
۱۹۶	چکیده
۱۹۷	آینده آرشیو وب
۱۹۷	اشاره
۱۹۷	خلاصه اجرایی
۱۹۷	بخش اول نمای کلی از چهار سناریوی احتمالی برای آینده:
۱۹۷	اشاره
۱۹۸	مقدمه
۲۰۱	ساختن آینده
۲۰۲	سناریوها [طرح ها]
۲۰۲	اشاره
۲۰۲	سناریوی نیروانا
۲۰۴	سناریوی آپوکالیپس (آخر الزمان)
۲۰۷	سناریوی انفرادی
۲۰۷	سناریوی غبار آلود
۲۱۱	مرور آینده
۲۱۱	یادگیری از وب پویا
۲۱۳	مصورسازی
۲۱۳	چالش ها
۲۱۵	برنامه های کاربردی جست و جو همانند شکارچی
۲۱۶	تحلیل های شبکه اجتماعی
۲۱۸	مثال ها:
۲۲۳	سنجش های دگرساز

۲۲۶	وب نوشت (حاشیه نگاری) اجتماعی
۲۲۸	معماران جدید
۲۳۱	ماشین های اجتماعی
۲۳۳	شبکه های نقشه برداری
۲۳۴	علم وب
۲۳۴	اشاره
۲۳۶	درک تجربه به جای محتوا
۲۳۷	تحلیل وب معنایی و مجموعه داده های پیوند شده
۲۴۰	چالش های فعلی و آینده
۲۴۰	اشاره
۲۴۲	وب مجتمع: زندگی آرشیو وب
	چالش بلند مدت: دو چالش در این سؤال نهفته است، که هر دوی آن ها مستلزم مشارکت کنندگان و فعالان زیادی خواهد بود. نخست: ما باید دوباره در مورد اینکه چگونه اینترنت را ببینیم و مهندسی کنیم فکر کنیم، در حال حرکت از موجودیتی تک لایه با پیوندهای جانبی
۲۴۴	وب در حال تغییر
۲۴۶	کاربرد های آرشیو ها و وب گاه ها
۲۴۹	متخصص وب
۲۵۰	وب دیداری
۲۵۵	وب همان گونه که بود
۲۵۸	ساختار وب
۲۵۸	ایده ها چگونه تکثیر می شوند
۲۵۹	وب غیر قانونی
۲۶۱	رد پای رقومی
۲۶۲	وب داده
۲۶۵	نتایج: مسیر پیش رو
۲۶۷	منابع
۲۷۳	فصل دوم: تجارب جهانی و مسائل بومی در آرشیو سازی وب
۲۷۳	اشاره
۲۷۴	چکیده
۲۷۵	آرشیو وب در دنیای وب ۲/۰ شعبه آرشیو وب و حفاظت رقومی کتابخانه ملی استرالیا
۲۷۵	اشاره
۲۷۵	مقدمه
۲۷۶	سه روش شناسی آرشیو
۲۸۲	جمع آوری فایل ها
۲۸۴	دستورالعمل های جمع آوری
۲۸۵	دستور عمل های آینده
۲۸۵	شرح حال مختصری از پدید آورنده
۲۸۶	چکیده
۲۸۷	آسیب شناسی زبان و خط فارسی در بازیابی اطلاعات: نگاهی به موتور های کاوش و پایگاه های برخط
۲۸۷	اشاره
۲۸۷	درآمدی بر مشخصه های زبان فارسی
۲۸۹	پیشینه پژوهش
۲۹۰	رسم الخط فارسی و بازیابی اطلاعات

۲۹۲	مسائل صرفی و بازیابی اطلاعات
۲۹۳	مسائل معنایی و بازیابی اطلاعات
۳۰۲	سخن پایانی: پیشنهادهایی در جهت بهبود وضعیت
۳۰۴	منابع
۳۰۸	چکیده
۳۰۹	ارزیابی کاربرد پذیری وبگاه نهاد کتابخانه های عمومی کشور
۳۰۹	اشاره
۳۰۹	مقدمه
۳۱۰	۲-روش پژوهش و توجیه رویی آن
۳۱۰	۳-شیوه گردآوری اطلاعات و تجزیه و تحلیل آن ها
۳۱۰	۴-تجزیه و تحلیل داده ها و ارائه یافته ها
۳۱۸	۵-نتیجه گیری
۳۱۸	۶-پیشنهادهای پژوهش
۳۱۸	منابع
۳۲۴	چکیده
۳۲۵	امکان سنجی برداش وبگاه ها در سازمان اسناد و کتابخانه ملی ایران
۳۲۵	اشاره
۳۲۵	مقدمه
۳۲۶	وبگاه ها: شناخته شده ترین منابع وب
۳۲۶	ارزیابی وبگاه ها
۳۲۶	استانداردها و طرحهای ابر داده ای سازماندهی منابع وب
۳۲۷	برداشت منابع الکترونیکی در سازمان اسناد و کتابخانه ملی ایران
۳۲۷	بیان مساله
۳۲۹	پیشینه پژوهش
۳۲۹	پیشینه پژوهش در ایران:
۳۲۹	پیشینه پژوهش در خارج از ایران:
۳۳۰	اهمیت پژوهش
۳۳۰	اهداف پژوهش
۳۳۲	بررسی های اساسی
۳۳۲	تعاریف عملیاتی
۳۳۲	جامعه آماری
۳۳۲	روش پژوهش و ابزار گردآوری داده ها
۳۳۲	یافته های پژوهش و پاسخ به بررسی های اساسی
۳۳۳	۱.استفاده از فهرست وب گاه ها
۳۳۳	۲.ورود اطلاعات وبگاه ها در وبگاه سازمان اسناد و کتابخانه ملی ایران
۳۴۰	بحث و نتیجه گیری
۳۴۲	پیشنهاد های برخاسته از پژوهش
۳۴۲	منابع
۳۴۶	چکیده
۳۴۷	مقدمه
۳۴۸	ایجاد آرشیو گزینشی منابع تحت وب

۳۴۸	اشاره
۳۴۹	تاریخچه
۳۵۱	آرشیو پاندورا در حال حاضر
۳۵۱	محدوده وظایف
۳۵۳	نیروی انسانی و افراد شاغل در پاندورا
۳۵۴	هزینه دستیابی به وب گاه ها و نشریات برخط
۳۵۴	محدوده برآورد هزینه
۳۵۶	روش شناسی انجام کار
۳۵۸	نحوه محاسبه هزینه ها
۳۵۸	هزینه های فراهم آوری منابع آرشیوی
۳۶۰	مقایسه با نوع جایی
۳۶۱	هزینه فعالیت های خاص
۳۶۱	امکان کاهش هزینه ها
۳۶۲	نتیجه گیری
۳۶۳	قدردانی
۳۶۳	منابع
۳۶۵	چکیده
۳۶۶	بایگانی وب علمی در مقیاس کوچک
۳۶۶	اشاره
۳۶۶	۱-جرائی بایگانی علمی در مقیاس کوچک
۳۶۷	۲-بایگانی دیجیتال برای مطالعات زبان چینی
۳۶۷	اشاره
۳۶۹	۲-۱-گام های اولیه
۳۷۰	۲-۲-توسعه پایدار سازمانی
۳۷۱	۲-۳-سخت افزار
۳۷۲	۲-۴-نرم افزار
۳۷۲	۲-۵-فرآیندها
۳۷۴	۲-۶-سیاست گذاری و خط مشی مجموعه
۳۷۶	۲-۷-مشارکت
۳۷۷	۳-درس های آموخته شده: جمع بندی
۳۷۷	۴-منابع مفید
۳۸۰	چکیده
۳۸۱	بررسی و مقایسه قابلیت های قالبهای یونی مارک و مارک ۲۱ برای سازماندهی منابع اطلاعاتی وب
۳۸۱	اشاره
۳۸۱	۱-مقدمه و بیان مسئله
۳۸۲	۲-بررسی های پژوهش
۳۸۳	۳-پیشینه پژوهش
۳۸۴	۴-روش شناسی پژوهش
۳۸۴	۵-تجزیه و تحلیل یافته ها
۳۸۶	۶-بحث و نتیجه گیری افزایش روزافزون منابع اطلاعاتی وب و آشننگی موجود در بازیابی اطلاعات جامع و مانع، لزوم سازماندهی منابع اطلاعاتی وب را بر رنگ تر می سازد از آن جایی که منابع وبی سازماندهی شده جزئی از نظام های ذخیره و بازیابی کتابخانه ای خواهند بود لازم است
۳۹۳	منابع

۳۹۶	چکیده
۳۹۷	سنجش رابط کاربر پایگاه های اطلاعاتی پیوسته: مجلات تمام متن فارسی
۳۹۷	اشاره
۳۹۷	مقدمه
۳۹۸	بیان مسأله و ضرورت پژوهش
۳۹۹	هدف پژوهش
۴۰۰	پیشینه ی پژوهش
۴۰۲	جمع بندی پیشینه
۴۰۳	روش شناسی پژوهش
۴۰۳	یافته های پژوهش
۴۱۰	بحث و نتیجه گیری
۴۱۱	پیشنهاد هایی برای پژوهش های آینده:
۴۱۱	منابع
۴۱۳	چکیده
۴۱۴	قانون واسپاری وب فرانسه: راهبردهایی برای گردآوری دامنه ملی
۴۱۴	اشاره
۴۱۴	۱. قلمرو فرانسه
۴۱۴	۱.۱. تعریف حوزه قانون و اسپاری
۴۱۹	۲.۱. در حال حاضر کجا هستیم
۴۲۱	۲. طراحی خزش
۴۲۱	۱.۲. هدف چیست؟
۴۲۳	۲.۲. فهرست هسته
۴۲۸	۳.۲. تنظیمات خزش گر
۴۲۹	۱.۳.۲. حوزه
۴۳۱	۲.۳.۲. اولویت های خزش
۴۳۳	۳.۳.۲. سایر تنظیمات
۴۳۵	۴.۲. پروتکل حذف رویت ها
۴۳۷	۵.۲. برنامه ریزی با همکاری آرشیو اینترنت (IA)
۴۳۸	۳. خزش گر در حال کار
۴۳۸	۱.۳. آزمایش خزش ها
۴۳۸	۲.۳. ارتباط کتابخانه ملی فرانسه و آرشیو اینترنت در طول خزش
۴۳۹	۳.۳. «خزش تکه ای»
۴۳۹	۴. پیامدهای خزش
۴۳۹	۱.۴. اشکال اصلی
۴۴۲	۲.۴. توزیع هر سرآیند
۴۴۶	۳.۴. فایل های ویدئویی
۴۴۷	۴.۴. توزیع هر TLD
۴۵۰	۵.۴. Robots.txt
۴۵۱	۶.۴. عمق خزش
۴۵۴	۷.۴. وب گاه های بزرگ
۴۵۴	۱.۷.۴. دامنه ها

۴۵۶	دانشنامه های سطح دوم
۴۵۷	۵ نتیجه گیری
۴۵۹	تشکر و قدردانی
۴۶۰	منابع
۴۶۵	چکیده
۴۶۶	معرفی آرشیوهای وب به عنوان یک خدمت جدید کتابخانه: تجربه کتابخانه ملی فرانسه
۴۶۶	اشاره
۴۶۷	چارچوب حقوقی
۴۶۹	دانشنامه
۴۷۱	ابزارها و روشهای جمع آوری
۴۷۲	کشف منبع: خدمات و ابزارهای دسترسی برای کاربران نهایی
۴۷۲	تعیین جای خدمت و محدودیت های دسترسی
۴۷۳	ابزارهای جستجو و دیدن
۴۷۸	مجموعه های ویژه
۴۸۲	کاربرد اطلاعات و ارقام
۴۸۲	تحلیل کمی
۴۸۴	تحلیل کیفی
۴۸۸	نظر سنجی ها
۴۹۱	راهبردهای به حساب آوردن آرشیوهای وب به عنوان بخشی از کار روزانه کتابخانه
۴۹۲	راهبردهایی برای رسیدن به کاربران نهایی
۴۹۴	منابع
۴۹۵	چکیده
۴۹۶	یک سال آرشیو وب گزینشی
۴۹۶	اشاره
۴۹۶	مقدمه
۴۹۷	۱.۲ آرشیو وب گزینشی در کتابخانه ملی زلاندنو
۴۹۷	۱-۲ انگیزه
۴۹۷	۲-۲ تاریخچه درو/هاروستینگ
۴۹۹	۳-۲ نرم افزار گردآوری وب
۵۰۰	شکل ۱. فهرست WCT
۵۰۲	۴.۲ اعضا و منابع
۵۰۴	۵-۲ سطح گردآوری / درو
۵۰۷	گردآوری با WCT ویرایش ۱.۱
۵۰۷	۱.۳ تجربه اولیه
۵۰۹	۳.۲ مشکلات گردآوری
۵۰۹	کمبودهای تورق نرم افزار
۵۱۱	نسخه های میانجی کاربر
۵۱۱	اشتباه «در مرحله ایست درگیر شدن»
۵۱۱	وب گاه های بزرگ
۵۱۱	۳.۳ خلاصه تجربه با ویرایش ۱.۱
۵۱۱	گردآوری با WCT با ویرایش ۱.۲

۵۱۱	۴.۱.نم افزار گردآوری تاریخ
۵۱۲	۲.۴.نم افزار تورق (مرورگر)
۵۱۲	۳.۴.نم افزار هرس درخت تصمیم
۵۱۲	۴.۴.وب گاه های بزرگتر جمع آوری می شوند
۵۱۳	۵.۴.گسترش ذخیره دیجیتال با ارزش
۵۱۳	۶.۴.ارتباط
۵۱۳	۷.۴.فهرست گردش کار و دسترسی
۵۱۴	۵.گردآوری رویداد انتخابات هیئت محلی
۵۱۶	۶.نتیجه گیری
۵۱۶	منابع
۵۱۸	درباره مرکز

سرشناسه: منتظر، غلامعلی، 1348 -، گردآورنده

عنوان و نام پدیدآور: مدیریت منابع اطلاعاتی وب [کتاب] / به کوشش غلامعلی منتظر و فرزانه شادان پور.

مشخصات نشر: تهران: سازمان اسناد و کتابخانه ملی ایران، 1391.

مشخصات ظاهری: ج2.

شابک: دوره 0-344-446-964-978 ؛ 150000 ریال: ج1 1-978-343-446-964-978 ؛ ج2 2-978-345-446-964-978 ؛

200000 ریال (ج2، چاپ اول)

وضعیت فهرست نویسی: فایا

یادداشت: ج2 (چاپ اول: 1391).

مندرجات: ج1. مبانی و تجربه های جهانی. -ج. 2. دیدگاه های فناورانه، اخلاقی و مدیریتی.

موضوع: وب--سایت ها--مدیریت

موضوع: منابع اطلاعاتی --مدیریت

موضوع: آرشوسازی وب

شناسه افزوده: شادان پور، فرزانه، 1344-، گردآورنده

شناسه افزوده: سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران

رده بندی کنگره: TK5105/888/م8م4 1391

رده بندی دیویی: 005/72

شماره کتابشناسی ملی: 3077380

دسترسی و محل الکترونیکی: <http://dl.nlai.ir/UI/2fb77759-f3eb-4f7f-ad9d-2cc0b917ed1d/Catalogue.aspx>

خیراندیش دیجیتال: انجمن مددکاری امام زمان (عج) اصفهان

ویراستار کتاب: خانم شهناز محققیان

ص: 1

اشاره

سازمان اسناد و کتابخانه ملی

جمهوری اسلامی ایران

مدیریت منابع اطلاعاتی وب

جلد اول

مبانی و تجربه های جهانی

به کوشش:

دکتر غلامعلی منتظر

و

فرزانه شادان پور

زمستان 1391

ص: 2

فهرست نویسی پیش از انتشار کتابخانه ملی جمهوری اسلامی ایران

سرشناسه: منتظر، غلامعلی 1348 - ، گردآورنده

عنوان و نام پدیدآور: مدیریت منابع اطلاعاتی وب / به کوشش غلامعلی منتظر و فرزانه شادان پور.

مشخصات نشر: تهران: سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران 1391

مشخصات ظاهری: 2 ج.

شابک: دوره: 0 - 344 - 446 - 964 - 978؛ ج. 1: 3 - 343 - 446 - 964 - 978؛

مندرجات: ج. 1. مبانی و تجربه های جهانی - ج. 2. دیدگاه های فناورانه، اخلاقی و مدیریتی

موضوع: وب--سایت ها -- مدیریت

موضوع: منابع اطلاعاتی-- مدیریت

موضوع: وب-- آرشیو سازی

موضوع: شادان پور، فرزانه، 1344- گردآورنده

شناسه افزوده: شادان، پور فرزانه 1344 - ، گردآورنده

شناسه افزوده: سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران

رده بندی کنگره: 00 5 /72

رده بندی دیویی: TK5105/88884 1391

شماره کتابشناسی ملی: 3077380

سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران

عنوان: مدیریت منابع اطلاعاتی وب، جلد اول: مبانی و تجربه های جهانی

به کوشش: دکتر غلامعلی منتظر (دانشیار دانشگاه تربیت مدرس) و فرزانه شادان پور (مربی، سازمان اسناد و کتابخانه ملی جمهوری

اسلامی ایران)

ویراستار ادبی: آرزو تجلی (کارشناس ارشد جامعه شناسی، سازمان اسناد کتابخانه ملی جمهوری اسلامی ایران)

تنظیم و تصحیح: مهشید برجیان (کارشناس ارشد کتابداری و اطلاع رسانی، سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران)

ویراستار استنادی (مقالات تألیفی): فروزان رضایی نیا کارشناس کتابداری و اطلاع رسانی، سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران)

نمونه خوانی و اصلاحات: مهشید برجیان، فاطمه رمضانپور، آمنه هزارخانی، زهرا زاهدی، محمد رضا میقانی، ملیحه حاجی زاده مقدم

طراحی جلد و صفحه آرایی: شهره خوری

ناظر فنی چاپ: نصرت الله امیرآبادی

ناشر: سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران

شمارگان: 500 نسخه

بها: 15000 تومان

نشانی: تهران بزرگراه شهید حقانی (غرب به شرق)،

بعد از ایستگاه مترو بلوار کتابخانه ملی

تلفن فروشگاه 81623318-81623315-88941946 دورنگار: 88947496

وب سایت: www.nlai.ir

پست الکترونیک انتشارات: Publication@nlai.ir

ص: 3

سازمان استاد و کتابخانه ملی جمهوری اسلامی ایران

مدیریت منابع اطلاعاتی وب

جلد اول

مبانی و تجربه های جهانی

ص: 4

سخن نخست ... نه

به جای مقدمه ... 1

فصل اول: مبانی مدیریت و آرشیو وب ... 3

بایگانی شبکه وب : مباحث و روش ها: نوشته ژولین ماسانه / ترجمه فهیمه باب الحوائجی ... 4

از آرشیو اینترنت تا آرشیو در اینترنت: نوشته میشل کیمتون / ترجمه مرضیه هدایت ... 62

کاربرد وب و مطالعات مربوط به آن: نوشته استیو جونز گمیل جانسون / ترجمه سید مهدی طاهری، سید محمد موسوی ... 74

خصوصیات وب ایران: نوشته مریم پیروزمند ... 90

آینده آرشیو وب: نوشته اریک تی. مه یر، آرتور توماس، رالف شرودر، مؤسسه اینترنت آکسفورد / ترجمه رضا خانیپور، محبوبه قربانی ...

106

فصل دوم: تجارب جهانی و مسائل بومی در آرشیو سازی وب ... 149

آرشیو وب در دنیای وب 2/0 شعبه آرشیو وب و حفاظت رقومی کتابخانه ملی استرالیا: نوشته ادگار کروک / ترجمه مرجان هادیزاده ...

150

ص: 5

آسیب شناسی زبان و خط فارسی در بازیابی اطلاعات: نگاهی به موتورهای کاوش و پایگاه های برخط: نوشته شعله، ارسطو پور فاطمه
احمدی نسب ... 158

ارزیابی کاربردپذیری وبگاه نهاد کتابخانه های عمومی کشور نوشته صدیقه محمد اسماعیل، ماهرخ ناصحی اسکویی ... 176

امکان سنجی پردازش وبگاه ها در سازمان اسناد و کتابخانه ملی ایران: نوشته رضا خانی پور، محبوبه قربانی، سهیلا فعال ... 190

ایجاد آرشیو گزینشی منابع تحت وب، بررسی هزینه های مربوط به فراهم آوری منابع تحت وب در کتابخانه ملی استرالیا: نوشته مارگارت
فیلیس / ترجمه صدیقه محمد اسماعیل ... 204

بایگانی وب علمی در مقیاس کوچک DACHS نوشته هانو لشر / ترجمه حمزه علی نورمحمدی ... 218

بررسی و مقایسه قابلیت های یونی مارک و مارک 12 برای سازماندهی منابع اطلاعاتی وب نوشته رقیه حجازی مهرداد کوبی ... 230

سنجش رابط کاربر پایگاه های اطلاعاتی پیوسته مجلات تمام متن فارسی نوشته صدیقه جعفرزاده، معصومه پیروزفر، عبدالحسین فرج
پهلوی ... 242

قانون و اسپاری وب فرانسه: راهبردهایی برای گردآوری دامنه ملی نوشته فرانس لاس، فارگوس کلمنت کیوری برت وندلاند / ترجمه سودابه
نوذری ... 254

معرفی آرشیوهای وب به عنوان یک خدمت جدید کتابخانه: تجربه کتابخانه ملی فرانسه نوشته سارا اویری / ترجمه زهرا تهوری ... 286

یک سال آرشیو وب گزینشی با WCT در کتابخانه ملی زلاندنو نوشته گوردون پنیتر، سوزانا جو، وائیتا لا لا، گیلیان لی / ترجمه احترام
السادات کیانمهر ... 304

از ویژگی های قرون گذشته بی خبری بود و تمایز جدی عصر جدید نسبت به گذشته دسترسی آسان به اطلاعات است. بشر با از سر گذراندن سه موج و پارادایم، کشاورزی صنعت و اطلاعات امروز در قرن بیست و یکم پا در عصر انفجار اطلاعات نهاده است این امر فی نفسه نه مطلوب است نه مذموم، بلکه به نحوه مدیریت ما نسبت به اطلاعات باز می گردد.

بشر امروزی به دلیل رشد روزافزون علم و فناوری در شرایط هشدارآمیز عدم قطعیت بسر می برد و همین مدیریت و تصمیم گیری را با چالش جدی روبرو ساخته است. اگر اطلاعات درست مدیریت شود و در تصمیم گیری ها به موقع به کار آید، و از دو ویژگی صحت و سرعت برخوردار باشد، می تواند منشأ تصمیم های تحول آفرین شود. ویژگی دیگر این عصر ظهور و حضور همه جانبه اطلاعات دیجیتال است. دورانی فرارسیده است که در آن بناست دانش مدون و تفکر مضبوط بشر علاوه بر کاغذ، و حتی بیش از آن، بر محمل «بیت» ها مسیر تولید، نشر و اشاعه، و مصرف را بیاماید. هم اطلاعات تولید شده تحت وب و هم میزان استفاده از این اطلاعات با سرعت فزاینده ای رو به رشد است. کشور ما بنابر اطلاعات وثیق از حیث تعداد کاربران و میزان حضور و فعالیت آن ها در وب جایگاه نخست را در منطقه خاور میانه داراست. این روند رو به رشد، با نصب العین قرار دادن آرمان های بلند انقلاب اسلامی در ترویج تفکر رهایی بخش اسلام ولایت مدار، وظیفه خطیری بر دوش نهادها و دستگاه های مسئول تولید، سیاستگذاری و نشر محتوا در محیط وب قرار می دهد و آن انجام بررسی های علمی و مستند به منظور ابتدای سیاستگذاری ها و عملکردها بر مبنای صحیح و کارآمد و متناسب با نیازهای گوناگون کاربران در این محیط است. اما وجه دیگر، صیانت از این محتوا و انتقال آن به نسل های آینده است که با توجه به ناپایداری محتوای قرار گرفته بر اینترنت و فناوری پیشرفته ای که برای چنین امر خطیری لازم است از اهمیت مضاعفی برخوردار می شود.

سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران بنابر مأموریت خویش دایر بر صیانت از میراث فکری کشور و اشاعه آن، عزم راسخ داشته است که برای مدیریت منابع اطلاعاتی مهم و رو به رشد وبی نیز چاره اندیشی نماید؛ بنابر این در سال 1389 نخستین بار در کشور به تهیه ساز و کار لازم برای ایجاد آرشیو ملی وب همت گماشته است.

از دیگر سو سازمان با علم به این که مدیریت در این حوزه مشارکت همه صاحبان اندیشه در حوزه تولید، سازماندهی و اشاعه اطلاعات تحت وب را می طلبد، مصمم شد «نخستین کنفرانس ملی مدیریت منابع اطلاعاتی وب» را برگزار نماید تا اهل علم و فناوری در این مجمع با هم اندیشی و تضارب آراء همچون

گذشته این سازمان را یار و یاور باشند.

این اثر مجموعه ای است فراهم آمده از تلاش پژوهشگرانی که با وجود نبودن مباحث مطرح شده در محورهای موضوعی کنفرانس، به ارائه ثمره پژوهش های خود همت نمودند؛ که با برگزیده ای از مقالات ترجمه ای در این عرصه پژوهشی ادغام و به طبع رسیده است رجاء واثق دارم که با الطاف الهی از این پس مدیریت منابع اطلاعاتی وب، و آرشیو وب به طور خاص، موضوع پژوهش و ابتکار عمل اهل دانش و فناوری در کشورمان قرار خواهد گرفت و در این عرصه نیز فرزندان این مرز و بوم تجسم گفتار نغز رسول اعظم صلی الله علیه و آله خواهند بود که «علم اگر تا ثریا، برود مردانی از فارس بدان دست خواهند یافت».

اسحق صلاحی

رئیس کنفرانس و رئیس سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران

ص: 8

گسترش روزافزون اطلاعات در شبکه اینترنت و سادگی بارگذاری انواع داده‌ها بر وب جهان را با شکل جدیدی از تولید، انتشار و مصرف اطلاعات مواجه کرده است. تغییر جایگاه شهروندان جامعه از مصرف‌کننده صرف اطلاعات به مولد و ناشر اطلاعات و فارغ از ساز و کارهای موسوم، سبب ساز روابطی جدید در عرصه ارتباطات اجتماعی و فرهنگ شده است. از سویی حجم رو به تزاید داده‌ها و چرخه عمر کوتاه اطلاعات موجود در وب، موجب شده که «گردآوری»، «پالایش»، «سازماندهی»، «ذخیره سازی» و «اشاعه» آن‌ها در زمره مسائل پژوهشی در نهادهای علمی و نیز بخش‌های پژوهش و نوآوری شرکت‌ها قرار گیرد؛ ضمن اینکه حفظ و دسترسی پایدار به اطلاعات موجود در وب، که خود جزئی از میراث فکری ملت‌ها محسوب می‌شود، به دغدغه‌ای جدی برای سازمان‌ها متولی حفظ و اشاعه میراث فکری به ویژه کتابخانه‌های ملی، بدل شده است.

این حوزه در جهان موضوعی نسبتاً جدید است و پیشینه آن به کمتر از پانزده سال می‌رسد، لیکن با سرعتی شتابان در حال رشد است و محققان مختلفی را از زوایای مختلف فنی، حقوقی، اقتصادی و حتی اخلاقی به سوی خود جذب کرده که گواه آن نیز طیف وسیعی از مقاله‌ها کتاب‌ها و گزارش‌های سازمانی است که در طی چند سال اخیر در سطح جهانی منتشر شده است. به رغم این نکات، در ایران همچنان این زمینه، حوزه‌ای بکر و کمتر مورد توجه محسوب می‌شود و در طی سالهای اخیر کمتر تحقیق بدان پرداخته، لیکن تزاید اطلاعات فارسی بر روی وب و برنامه‌های ملی کشور مبنی بر توسعه کاربری‌های مختلف بر شبکه‌های اطلاعاتی (از جمله توسعه دولت الکترونیکی، یادگیری الکترونیکی و کتابخانه‌های دیجیتالی) لزوم توجه به این موضوع را بیش از پیش نمایان می‌سازد. به همین دلیل سازمان اسناد و کتابخانه ملی جمهوری اسلامی همزمان با برگزاری «نخستین کنفرانس ملی مدیریت منابع اطلاعاتی وب» در صدد برآمد تا این حوزه را هرچه بیشتر به متخصصان و پژوهشگران باشناساند. کتاب پیش رو حاصل همین نیت متولیان این موضوع مهم است.

این کتاب مجموعه‌ای قریب به 30 مقاله برگزیده از مهم‌ترین منابع علمی منتشر شده در جهان و نیز قریب به 15 مقاله برگزیده از صاحب نظران ایران است که در قالب دو جلد تقدیم

حضور خوانندگان ارجمند می شود. این مقالات در چهار موضوع اصلی به شرح زیر تقسیم شده اند:

• مبانی مدیریت و آرشیو وب

• تجارب جهانی و مسائل بومی در مدیریت و آرشیو وب

• مسائل فناوریانه

• مسائل اخلاقی و مدیریتی

بی گمان این مجموعه می توانست به افزودنی های دیگر (هم از منابع خارجی و هم از دیدگاه سایر متخصصان ایرانی) به اثری پربارتر بدل گردد لیک نخستین گامی است که در این حوزه برداشته شده و مطمئناً در مراحل بعدی با همت سایر اندیشمندان، ویراست هایی غنی تر از آن حاصل خواهد آمد. نگارنده امیدوار است این مجموعه به مثابه بذری باشد که در کشتزار ذهن پژوهشگران کاشته شده و ان شاء الله در آینده ای نه چندان دور به نهالی پر طراوت در عرصه علم و عمل در جامعه اسلامی مان مبدل گردد.

در پیدایی این اثر کسان بسیاری همراهی و همکاری داشته اند که مقدم بر همه اندیشمندانی است که متن هر مقاله به خامه دانش افزای آنان امکان وجود یافته است. از این رو نگارنده سپاس فروتنانه خود را نثار نگارندگان و مترجمان ارجمند این اثر می نماید. گردآوری، تنظیم و آماده سازی مطالب کتاب به همت خانم ها فرزانه شادان پور و مهشید برجیان بوده و ویراستاری آن را خانم آرزو تجلی بر عهده داشته اند. ویراستار استنادی مقالات تألیفی را سرکار خانم فروزان رضایی نیا به انجام رسانده اند و سرکار خانم دکتر میترا صمیعی زحمت چکیده نویسی شماری از مقالات را که فاقد چکیده بودند متقبل شدند. نمونه خوانی و اصلاحات اثر حاصل تلاش خانم ها مهشید برجیان فاطمه رمضانپور آهنگری، آمنه هزار خوانی، زهرا زاهدی، ملیحه حاجی زاده مقدم و آقای محمد رضا میقانی بوده است. ضمن اینکه زیبایی متن و صفحه آرایی آن مدیون حسن سلیقه سرکار خانم شهره خوری است. زحمات لیتوگرافی، چاپ و صحافی کتاب نیز بر عهده جناب آقای امیر آبادی بوده که بر خود فرض می داند از همه این بزرگواران صمیمانه تشکر کند.

بی گمان پدید آمدن این اثر به همت مسؤولان گرانمایه سازمان اسناد و کتابخانه ملی جمهوری اسلامی بوده است و نگارنده امیدوار است خداوند آنان را در مسیر خدمت به فرهنگ و دانش ایران اسلامی مورد تأیید قرار دهد.

اللهم وفقنا لما تحب وترضی

غلامعلی منتظر

تهران- بهمن ماه یک هزار و سیصد و نود و یک خورشیدی

ص: 2

فصل اول: مبانی مدیریت و آرشيو وب

اشاره

ص: 3

بسیاری از جنبه های اجتماعی، اتفاقی هستند یا به طور کلی درباره اینترنت و به ویژه درباره وب بازتاب یافته اند. محافظت از وب، به این دلیل ضرورتی فرهنگی و تاریخی است. اما وب، از نظام های انتشاراتی قبل نیز برای ضرورت یک بازبینی ریشه ای از عملکردهای حفاظتی مرسوم، متفاوت می باشد. مفهوم میراث جمعی مشترک، شامل هر مصنوع محصول انسانی می شود از بناهای تاریخی معماری گرفته تا کتاب های جدید قرن بیستم ولو اینکه با فعالیت هایی حفاظتی مرتبط باشند (مثل آن ها که به طور اصولی و اختیاری سازماندهی شده اند) و قبلاً نمایان شده اند. در عصر حاضر علت بایگانی، دلیلی است که بر جزئیات تأکید میکند، به گونه ای پایدار ما را به اجتناب از عمومیت گرایی هدایت میکند و به طور خاص و منحصر به فرد که حوادث و رویدادها را مورد بررسی قرار می دهد. در حقیقت، امکاناتی که وب برای انتشار ایجاد می کند، منبعی منحصر به فرد از محتوا را ارائه می دهد که استدلال بایگانی وب، معقولانه انجام و تصدیق شده است. با توسعه شبکه های وب، به طور چشمگیری، حجم آن چه که می تواند منتشر شود و همچنین تعداد «ناشران» بالقوه یا خالقان محتوا با تقلیل هزینه های انتشار به تقریباً هیچ افزایش یافته است. دلگرم کننده است که ببینیم که بسیاری از مؤسسه های (حفظ) میراث، در بایگانی وب در حال به کارگیری هستند. بررسی اخیر توسط گروه پژوهش کتابخانه (RLG 2006) نشان داد که 60 درصد اعضای مورد بررسی اشان، بایگانی وب را قسمتی از مأموریت خود پنداشته اند.

* بایگانی شبکه وب: مباحث و روش ها (1)

ژولین ماسانه (2) ترجمه: فهیمه باب الحوائجی (3)

1- مقدمه

محصولات فرهنگی گذشته همیشه نقش مهمی در اطلاع رسانی و خود ادراکی جامعه و ساختن آینده آن ایفا می نمایند. شبکه جهان گستر وب (به طور خلاصه وب) رسانه ای جامع و گذراست در جایی که با درکی عمیق، فرهنگ مدرن در جهت یافتن شکل طبیعی عبارات و اصطلاحات است. انتشار، مباحثه، ایجاد، کار و تبادل اجتماعی در یک درک عمیق: بسیاری از جنبه های اجتماعی، انتقادی هستند یا به طور کلی درباره اینترنت و به ویژه درباره وب بازتاب یافته اند. محافظت از وب، به این دلیل ضرورتی فرهنگی و تاریخی است. اما وب، از نظام های انتشاراتی قبل نیز برای ضرورت یک بازبینی ریشه ای از عملکردهای حفاظتی مرسوم، متفاوت می باشد.

این فصل، بررسی موضوع های درخواستی از حفاظت از وب را ارائه می دهد و روش هایی که تا این تاریخ برای غلبه یافتن بر آن ها توسعه یافته است. ابتدا استدلال هایی را در برابر ضرورت و احتمالات بایگانی وب مطرح می کنیم. سپس، سعی می کنیم تفاوت های بسیار برجسته ای که وب را از دیگر

ص: 5

Web Archiving: Issues and Methods: in Masanes, Julien (ed.), Web Archiving. Berlin Heidelberg New - 1
York: Springer.pp.1-46

Julien Masané s -2

3- دانشیار کتابداری و اطلاع رسانی و دانشیار گروه کتابداری و اطلاع رسانی دانشگاه آزاد اسلامی واحد علوم و تحقیقات تهران

محصولات فرهنگی متمایز می کند، ارائه دهیم و دلالت های شان را برای محافظت ترسیم می نماییم.

این اساس وب را در بر می گیرد و آن را به عنوان نوعی سیستم نشر فعال و یک ابررسانه ویرایش اجتماعی و با یک محصول فرهنگی سراسری مورد بررسی قرار می گیرد. احتمالات و محدودیت های حفاظت هر یک از این جنبه های وب را مطرح می کنیم. سپس، رویکردهای روش های بررسی مهم را برای اکتساب سازماندهی و ذخیره سازی محتویات وب ارائه می کنیم.

فصل 2 و 4 و 5، جزئیات بیشتری درباره روش های بررسی و ابزارهایی برای اکتساب مضامین تهیه کند و فصل های 6-8 بر روی دسترسی استخراج اطلاعات و حفاظت از محتوای وب تمرکز کرده است. دو فصل آخر این کتاب بررسی هایی موردی را ارائه می کند. بایگانی اینترنتی که بزرگ ترین بایگانی وب در جهان است (فصل 9) و DACHS یک پژوهش که در جهت گزینش بایگانی وب (فصل 10) می باشد. این فصل مقدمه ای کلی برای این کتاب مورد بررسی قرار گیرد:

در نهایت، موارد اولیه را در این حوزه فراهم می کند و طبقه بندی بایگانی های وب را برای ترسیم وضعیت جاری حفاظت وب را پیشنهاد می نماید.

- میراث جامعه و وب

2-محافظت از میراث

مفهوم میراث جمعی مشترک، شامل هر مصنوع محصول انسانی می شود، از بناهای تاریخی معماری گرفته تا کتاب های جدید قرن بیستم، ولو اینکه با فعالیت هایی حفاظتی مرتبط باشند (مثل آن ها که به طور اصولی و اختیاری سازماندهی شده اند) و قبلاً نمایان شده اند. شکل، اهداف، و کارایی حفاظت از میراث، به طور قابل توجهی با زمان و رسانه های بررسی شده متفاوت اند و تلاش این مقاله برای خلاصه کردن این تکامل کافی نیست. اجازه دهید فقط به یاد آوریم که از آماده سازی روشنفکری مذهبی (طبق کتابخانه کاسیودوروس و یواریوم، ریچه 1998) (1) برای ایجاد مجموعه ای به عنوان نشانه های قدرت (ابتکارات موزه مدرن توسط مدیسیس (2) در فلورانس در اواخر قرن پانزدهم را ببینید) برای کنترل وضعیت نظام مند حفاظت از فرهنگ ملی (ابتکار سپرده قانونی فرانسیلر (3) را ببینید) استفاده شده است و از انگیزه های گوناگونی در جهت جمع آوری نظام مند و حفاظت از محصولات فرهنگی در تاریخ، ناشی شده است.

در عصر حاضر، بایگانی ها به طور کلی، تمایل زیادی به فراگیر شدن دارند (آسبورن 1999) (4). همان طور که مایک فیدراستون (5) اظهار می کند:

علت بایگانی، دلیلی است که بر جزئیات تأکید می کند، به گونه ای پایدار ما را به اجتناب از عمومیت گرای هدایت می کند و به طور خاص و منحصر به فرد که حوادث و رویدادها را مورد بررسی قرار می دهد. تمرکز بر بایگانی مؤثر کالاها به طور فزاینده ای تغییر یافته است به سوی جزئیاتی با واقعیتی

ص: 6

Franceisler -3

osborn -4

Mike Featherstone -5

عظیم از زندگی روزمره دنیوی تمرکز یافته است (فیدرستون 2000) (1).

در حقیقت، امکاناتی که وب برای انتشار ایجاد می کند، منبعی منحصر به فرد از محتوا را ارائه می دهد که استدلال بایگانی تمایل به ستایش آن دارد. از این رو، می توانیم فرض کنیم که قانونی بودن بایگانی وب، معقولانه انجام و تصدیق شده است. با وجود این، محافظت از وب، مورد سؤال قرار گرفته و تاکنون مورد قبول همه واقع نشده است. استدلال ها در برابر بایگانی وب می تواند در سه مقوله دسته بندی شود. آن ها که مبتنی بر محتوای یافت شده در وب هستند، آن ها که تصور می کنند که وب خود - محافظ است، و آن ها که فرض می کنند بایگانی وب امکان پذیر نمی باشد.

2-1- به اندازه کافی خوب نیست؟

نخستین مقوله، بحث هایی را درباره کیفیت محتوای وب در بر دارد که گفته می شود فرض بر عدم مطابقت معیارهای مورد نیاز برای محافظت می باشد. این موقعیت، مدت زیادی توسط برخی متخصصان نشر جهانی (ناشران و کتابداران) حفظ شده است و در جهت تهدید بزرگی که توسط این رسانه جدید برای بقا اعمال می شود همراه می شود. معمولاً با تأکیدی درباره میزان گسترده اطلاعات وب و فقدان دانش درباره روش های بایگانی وب و هزینه های آن همراه شده است.

منافع این موقعیت، از انتقال سیستم نشر پیوسته و آن هایی که بر حفاظت و بازدهی صنعت نشر پیوسته ادامه می دهند آگاهی دارد. اما آن ها از مرزهای آن چه که محافظت شده است را اجتناب می ورزند همان اندازه که شبکه وب، محدودیت آن چه را که انتشار یافته است را گسترش می دهد. معادلات اقتصادی تولیدات فیزیکی حامل دانش (ادواری ها، کتاب، و مانند آن) میراث انقلاب گوتنبرگ هستند که باید طبق این دیدگاه بر استقرار محدودیت ها برای کاری که باید حفاظت شود ادامه دهند. حتی زمانی که این معادلات عمیقاً تعدیل شده باشند. از نظر تاریخی، این حقیقت که چه چیزی می تواند منتشر شود، به واسطه هزینه های فیزیکی محدود می شود (که شامل تولید و نقل و انتقال، ذخیره سازی، و هزینه های اداره است) و فیلترسازی را ایجاد می کند، برای آن که نظام های نشر بیشتر از پنج قرن است که این کار را انجام می دهند. اما این مورد مناسبی نیست و نسبتاً میراث تعادل ثابتی از قرن پانزدهم است که نقض شده است. توسعه شبکه های وب، به طور چشمگیری، حجم آن چه که می تواند منتشر شود و همچنین تعداد «ناشران» بالقوه یا خالقان محتوا با تقلیل هزینه های انتشار به تقریباً هیچ افزایش یافته است. مباحث درباره ارزیابی کیفیت، به ناچار ذهنی، در واقع با پنهان نمودن مذاکرات واقعی درباره گسترش جو انتشار می یابد.

اگر چه رشد ادواری ها در پایان قرن نوزدهم، قابل مقایسه با وسعت آن در انقلاب جاری نمی باشد، برخی خصوصیات (مثل در همکرد نوع نشر با وضعیت جسمانی و فکری) را به اشتراک گذاشته و عکس العمل های مشابهی را به وجود آورده است. گاهی اوقات برای جامعه کتابخانه برای مثال برای پذیرش این نوع انتشار در قفسه هایشان و همچنین در قلب شان، در بر گرفته شده است. همان طور که

ص: 7

فایت - شایب (2000) (1) برای موردی در فرانسه نشان داده است، رفتار توصیفی خاص که در سطح عنوان به آن نیاز داشت، توسط این جامعه مورد اغماض قرار گرفت و بخش کاملاً جدید مدیریت اطلاعات در کنار کتابخانه هارا، به همین دلیل، به وجود آورد (مستندسازی، نمایه سازی، نوشته های علمی). مباحث بر روی بایگانی وب، برخی شباهت هایی را بر حسب این رویدادهای فرعی، به اشتراک می گذارد که اگر در همان حالت باقی بماند، دیده خواهد شد.

فیلترسازی، اگر چه به مدت طولانی نیازی به تخصیص تولیدات فیزیکی منابع کار آمد دانش را ندارد، به طور کامل از بین نرفته است. بلکه از نقشی مرکزی به یک نقشی پیرامونی تغییر می یابد و باز هم در برخی فضاها مورد نیاز است (برای مثال، اعتبار سنجی علمی) و به شکل های جدید تجربه می شود (مثل ویکی پدیا، اسلش دات، بوگس فر (2))

چنان چه اکسل برونز (3) توضیح می دهد:

عکس العمل ناگهانی این فعل و انفعالات و وسایل ارتباط جمعی مشارکتی بسیار محسوس است. اگر کسی بتواند یا حداقل توانایی بالقوه آن را، [مثل یک] یک ناشر دارد، چه اثراتی بر روی مؤسسه های انتشاراتی موجود خواهد داشت؟ اگر اطلاعات موجود در وب بتواند با اطلاعات وسیع متنوع به آسانی ارتباط (4) داشته باشند، چه اثراتی بر روی چارچوب های (5) انتشاراتی سنتی می تواند داشته باشند؟ اگر توانایی بالقوه برای مخاطبان وب برای مشارکت در تولید و ارزیابی تعاملی محتوا وجود داشته باشد، برای ایجاد نقش های تولید کننده و مصرف کننده در وسایل ارتباط جمعی چه اتفاقی می افتد؟ (برانز، 2005).

در زمینه حفاظت این مسئله به طور جدی مورد توجه قرار گرفته است. یک موضوع قطعی است و آن مدینه فاضله ای است که امیدوار باشیم که تعداد کمی از کتابداران فیلتر سازی ناشران را در مقیاس وب جهانی جایگزین نمایند حتی اگر آن ها سنت مدیدی در انتخاب محتوا داشته باشند، این کار را در محیطی ساختار یافته تر انجام می دهند که چندین مرتبه در ابعاد کوچک تر داشته باشد. اگر چه این کار هنوز احتمالی است و برای جامعه ای که خوب تعریف شده و اهداف کوچک مفید است (فصل درباره انتخاب روش شناختی ها و فصل 10 را درباره داجز (6) پژوهشی برگرفته درباره بایگانی وب، پژوهش انجام شده و هم چنین بروگر (7) 2005) را ببینید) به کارگیری این موارد به عنوان یک ساز و کار جهانی برای بایگانی، وب واقع گرایانه نیست اما این حقیقت که گزینش دستی محتوا برای وسعت وب مقیاس گذاری نشده است دلیلی برای نپذیرفتن بایگانی وب نیست این فقط دلیل خوبی برای بررسی مجدد موضوع گزینش و کیفیت در این محیط است.

آیا می تواند بر اساس یک ارزیابی کیفی توزیعی در سطح بالا و جامع باشد؟ این ارزیابی، به طور

ص: 8

Fayet-schibe -1

impact bogosphere, Slashdo -2

Axel Bruns -3

link -4

format -5

DACHS -6

ضمنی در دو سطح ساخته می شود:

کاربران وب به وسیله دسترسی بر محتوا، خالقان (محتوا) با پیوند دادن شکل محتوای صفحه های شان (ما در اینجا به قضاوتی که توسط خود خالقان قبل از اینکه محتوای شان را به صورت برخط در وب بگذارند در نظر نمی گیریم، که اگر به عنوان یک معیار انتخاب استفاده شود، به معنای بایگانی هر چیز است) همچنین، می تواند به طور ضمنی توسط افزایش انتخاب کنندگان فعال ایجاد گردد.

اجازه دهید ابتدا دسترسی کاربران را بررسی کنیم. گسترش جو انتشار پیوسته تحت چیزی که ظرفیت اقتصادی برای چاپ فیزیکی مجاز دانسته است، نتایج دیگری را در بر می گیرد: افت مکانیکی در میانگین تعداد خوانندگان هر واحد از محتوای منتشر شده برخی صفحه ها حتی نه تنها توسط هیچ انسانی خوانده نمی شوند بلکه توسط هیچ رباتی هم نمایه نمی شود. بوفخاد و ویونات (2003) (1)، استفاده از سیاهه های مربوط و پرونده های سرور یک وبگاه بزرگ دانشگاهی را نشان داده اند که 5 درصد صفحه ها فقط توسط ربات ها قابل دسترسی بودند و به 25 درصد آن ها هرگز دستیابی نداشتند، بدان معنا که توزیع دسترسی به محتوای پیوسته بسیار طولانی مدت ارائه خواهد شد.

اما این تکامل در نشر مدرن، کاملاً جدید نیست. رشد و درجه بالایی از تخصصی شدن انتشارات ادواری تقریباً الگوی مشابه دسترسی را نشان می دهد. آیا این استدلالی برای عدم حفظ ادواری هاست؟ در بیش تر کشور ها، نظام های سپرده قانونی، به طور مستقل، از انتشارات آن چه را که مورد استفاده قرار می گیرند، را نگهداری می کنند. این بی تکلیفی علائق خوانندگان آینده پیش بینی می کند.

مطمئناً برای حفاظت کنندگان در ارزیابی مفید بودن محتوای پیوسته برای نمایش و تلاش برای پیش بینی برای آینده، تا زمانی که برای جوامع کاربران از پیش تعریف شده باشد، امکان پذیر است.

الگوهای دسترسی همچنین می توانند برای اداره نظام های بایگانی جهانی استفاده شوند: در مورد بایگانی وب اصلی، تاکنون، مجموعه بایگانی اینترنت توسط الکسا (2) اهدا شده است که از الگوی دستیابی برای تعیین عمق خزش برای هر سایت استفاده کرده است (فصل 9 کیمپتون) (3) و دیگران (2006) را ببینید). همچنین، می تواند به وسیله پرس و جو های فرستاده شده برای موتور جست و جو اجرا گردد (پانندی و اولتسون، 2005) (4). اما سؤال کلیدی برای بایگانی های وب این است: چگونه به این اطلاعات دست یابیم و کدام مرز مورد استفاده قرار می گیرد؟

ترافیک اطلاعات، معمولاً وجود ندارد و موتورهای جست و جو از ابتکار الکسا پیروی می کنند و آن را از میلیون ها نوار ابزار نصب شده در مرورگرهایی که از اطلاعات شناوری (5) کاربران به آن ها می گذرند به دست می آورد.

مؤسسه های بایگانی از کجا می توانند به آن [اطلاعات] دسترسی پیدا کنند چنان چه آن ها خودشان جست و جوی ویژه ای را ارائه ندهند؟ مرزهای آن چه باید باشند؟ آیا باید در صفحه یا در سطح سایت

ص: 9

Kimpton et al -3

Pandey and Olston -4

navigation -5

به کار برده شود؟ (الکسا در سطح سایت از آن استفاده می کند) آیا عمق خزش را فقط در سطح نخست هر سایت محدود می کند (که این بدان معناست که حداقل در سطح اول هر سایت در تمام موارد اشغال خواهد بود)؟

حتی اگر این معیار، مباحث اجرایی عملی زیادی را نمودار سازد، مزیت بردن به عنوان حمل کننده برای تمرکز بر بایگانی، ستانده های میلیون ها کاربر - نه اجتماعی کوچک - را دارد که به خوبی با مدل انتشار جامعه وب تطبیق یافته است.

معیار دیگر، سطح اهمیتی است که توسط درجه پیوند درونی یک صفحه (یا یک سایت) اندازه گیری می شود. این مسئله استدلال شده است (میسائز، 2002) (1) که این تعادل مناسب در محیطی فرامتن، با درجه ای از عمومیت است که ویژگی های انتشارات سنتی را مشخص می کند و عملاً مزیتش در قابلیت استفاده برای استخراج ماتریس های پیوندی وب است (پیچ) (2) و همکاران،

1998؛ ایتبول (3) و همکاران، 2003، 2002؛ پاستور - ساتوراس و وسپینگانی (4) (2004). روش دیگری در انبوه کردن ارزیابی کیفی ایجاد شده است، البته نه توسط کاربران، بلکه به وسیله ایجاد کنندگان صفحه (و پیوندها) این مدل ارزیابی کیفی توزیع شده به خوبی با طبیعت انتشار بر روی اینترنت و به طور عملی با امکان اجرا مطابقت دارد.

سرانجام، ممکن است توسط مشارکت های بیشتر، وظیفه گزینش مطالب برای بایگانی کردن مقیاس را بالا برد. این کار می تواند با در بر گیری مؤسسه های بیشتری در انجام و تسهیل توسط ایجاد خدمات بایگانی انجام شود که با بخش فنی مسئله سروکار دارد. این کار به وسیله بایگانی خدمات در بایگانی اینترنتی پیشنهاد شده است که در سال 2006 شروع شده و عبارت است از توانمند سازی تنظیم و مدیریت آسان برای کتابخانه و آرشیو که نمی تواند در زیر ساختهای عملکردی مورد نیاز برای بایگانی وب سرمایه گذاری کند.

تحول ممکن دیگر با عمومیت بخشیدن به این توانایی برای هر کاربر وب در مشارکت فعال می باشد؛ البته چنانچه بخواهند در بایگانی وب شرکت کنند. انگیزه اصلی کاربران در این مورد، سازماندهی حافظه وب شخصی برای امکان بازگشت به مرجع بعدی برای محتوای با ثبات و، کاوش و سازماندهی آن به عنوان راهی برای مبارزه با علائم «گم شدن در فضای مجازی» است چندین بررسی در مورد کاربران نشان داده است که حفظ نشانه های محتوای مشاهده شده برای بسیاری از کاربران ضروری است (تیوان 2004) (5) البته آن ها از روش های مناسبی نیز استفاده می کنند (6) (جونز 2003، 2001). در بایگانی وب شخصی، دنبال کردن پیشینه کاربر بر روی وب می تواند برای یک سازمان دهی شخصی و محوریت زمانی

ص: 10

Masanes -1

Page -2

Abiteboul -3

Pastor-Satorras and Vespignani -4

Teevan -5

Jones -6

در حافظه وب میسر شود (رکیموتو 1999) (1) دامیس و همکاران (2003) (2)، رینگر و همکاران، (2003) (3). خدمات پیوسته متعددی (Furl, My Yahoo) قبلاً در بایگانی وب شخصی در سطح صفحه ارائه شده که با قابلیت علامت زدن ترکیب شده است. خدمات بایگانی های هانزو (4)، حوزه گسترش یافته (مضمون، تمام سایت) و همچنین در هم و بر همی قابلیت های بایگانی، با ابزارها و خدمات دیگر (مثل بلاگ ها، مرورگرها، و مانند آن) از طریق باز کردن ای پی آی (5) را اجازه می دهد که آن با یک بایگانی سرویس گیرنده با قابلیت های P2P توسعه بیشتری خواهد یافت؛ و به طور چشمگیری امکانات برای کاربران در ثبت تجربیات شان از وب را به عنوان بخشی از زندگی رقومی توسعه خواهد داد (فریمن و گلنتر 1996) (6)، گمل و همکاران، 2002). در مورد استفاده بالقوه از مخازن کاربران در یک بایگانی وب نظیر به نظیر نیز می توانید مانتراتزیس آرگون (2004) (7) را ببینید.

دیده شده است اگر این گسترش و دموکراسی شدن نقش بایگانی بتواند مانند تفسیر و سازماندهی اطلاعات توسعه یابد موجب پیشرفت علامت گذاری (8) گولدر و هامبرمن 2005 (9) و نظام های بلاگ کردن می شود (هالاویس 2004) (10)، ؛ برونز 2005 (11). که در این صورت کمک ارزشمندی خواهد بود و برای محافظت سازمان هایی که می توانند نظارت طولانی مدت بر این محتواها داشته باشند، درونداد [خوبی] خواهد بود.

همان طور که دیدیم استدلال ها در برابر بایگانی وب بر اساس کیفیتی است که درباره فرضیاتی بر پا شده است مانند (1) کیفیت محتوا تحت فضای سنتی محتوای ویرایش شده به طور مرسوم، کافی و مناسب نیست؛ و (2) فقط گزینش دستی و تک به تک که توسط حفاظت کنندگان ایجاد شده، می تواند جایگزین فقدان فیلترسازی ناشران شود (روشی که نمی تواند درست به اندازه مقیاس وب باشد، در حالی که همه با (فیلپز 2005) (12)، موافق اند. این دو استدلال، فقدان درک اساس توزیعی وب را نشان می دهند و اینکه چگونه می تواند وسیله نفوذ سازماندهی حافظه اش در مقیاس بزرگ باشد.

2-2-1- خود بقایی رسانه؟

دومین مقوله استدلال ها اظهار می دارد که وب رسانه ای خود بقاست. در این دیدگاه، منابعی که برای حفاظت شدن مناسب هستند، بر روی سرورها نگهداری می شوند. بقیه به اراده به وجود آورنده اصلی، ناپدید خواهد شد. از آن جا که نوع اول استدلال درباره کیفیت تقریباً در مجموعه برنامه های جهانی

ص: 11

Rekimoto -1

Dumais et al -2

Ringel et al -3

Hanzo -4

API -5

Freeman and Gelernter -6

Orgun, Mantratzis -7

Tagging -8

Golder and Huberman -9

Halavais -10

Bruns -11

Phillips -12

یافت شده است، این‌ها بیشترین طرفداران را در علوم کامپیوتر جهانی به دست آورده‌اند. اگر چه به شدت در روزهای اول پشتیبانی شده است، باید بگوییم که همانطور که زمان می‌گذرد و محتوا از وب محو می‌شود این مسئله کمتر مورد [چالش] است. اسناد مطالعاتی زیادی بر طبیعت بی‌دوام منابع وب، این ادعا را که که وب رسانه‌ای خود بقاست دچار شکست می‌کند. برای مثال، برای مروری بر ادبیات موضوع (کهلر 2004) (1) و (اسپاینلیس 2003) (2) را ببینید. این مطالعات بر دسترس پذیری منابع با همان URL تأکید می‌کنند و نه تغییرات بالقوه‌ای که می‌تواند متحمل شود. مطالعات نشان می‌دهند که میانگین نیم عمر هر صفحه وب (مدت زمانی که نیمی از صفحه‌ها ناپدید خواهد شد)، فقط دو سال است. این بررسی‌ها بر روی دسترس پذیری منابع در URL مشابه متمرکز هستند و تغییرات بالقوه‌ای وجود ندارد که آن‌ها متحملش شوند. همچنین برخی، مضمون‌ها را مورد تحقیق و بررسی قرار داده و میزان تغییر را ارزیابی می‌کنند. چو و گارسیا - مولینا (2000) (3)، نیم عمره 5 روزه‌ای را برای میانگین صفحه‌های وب کشف کردند. فترلی و همکارانش (2003) (4) نشان دادند چگونه این میزان تغییر با اندازه و موقعیت محتوا مرتبط می‌باشد.

دلایل بسیار زیادی برای تمایل منابع به ناپدید شدن از وب وجود دارد. نخست، محدودیت زمانی مجاز شدن دامنه نام (معمولاً 1-3 سال) است که به وسیله طراحی هر فضای وب در یک انتقال و شرایط غیر ایمن واقع می‌شود.

دوم، توان الکتریکی پایدار، پهنای باند و سرورهایی که به پشتیبانی انتشارات نیاز دارند - همان طور که در مقابل ماهیت خارجی انتشار صورت می‌گیرد. اما حتی وقتی که نامیدن فضا و منابع نشر، مصون هستند، سازماندهی و طراحی اطلاعات می‌توانند نقش مهمی را در حالت ارتجاعی منابع روی سرورها ایفا کند (برنرز لی 1998) (5). همان طور که برنرز مخترع وب ادعا کرده است:

«اصلاً استدلالی در نظریه برای افراد وجود ندارد تا URL‌ها را تغییر دهند (یا اسناد نگهداری شده را متوقف کنند) اما میلیون‌ها استدلال در عمل وجود دارد (برنرز، 1998).

تغییر افراد، سازماندهی داخلی، طرح‌ها، فناوری‌های سرور وب، عملیات نام دادن، و مانند آن می‌تواند ناشی از بازسازی و گاهی فقدان اطلاعات باشد.

از این دیدگاه، سبک رشد نظام مدیریت محتوا (6) در انتشار، برداشت‌های گمراه‌کننده در برقراری نظم هنگام بحران - چنانچه نظام مدیریت محتوا آورد چون معمولاً یک سبک ساختاری از اطلاعات یکپارچه و اغلب قابلیت‌های بایگانی کردن را دارا می‌باشد. مسئله این است که آن‌ها به لایه‌های دیگر وابستگی به نرم افزار اضافه می‌شوند (نرم افزار نظام مدیریت محتوا) چون استاندارد سازی در این حوزه وجود ندارد. معماری‌های اطلاعات بر اساس نظام مدیریت محتوا ثابت شده است که خنک است تا زمانی که نظام

ص: 12

Koehler -1

Spinellis -2

Cho and Garcia-Molina -3

Fetterly etal -4

Berners-Lee -5

CMS -6

مدیریت محتوا تغییر نکند، یعنی که خیلی طولانی نباشد.

اما آیا طراحی اطلاعات دستی است یا توسط سیستم انجام می شود. وب رسانه ای خود بقا نیست و نخواهد شد. مهم ترین دلیل آن مغایرت فعالیت های انتشار و حفاظت می باشد. انتشار، به معنای ایجاد تازگی است حتی زمانی که در مطالب کهنه هزینه شده باشد و (برای مثال در یک فضای نام گذاری مشابه یا کتاب های جدید و قدیمی باید در مخزن ناشران مشابه، با هم قرار گیرند).

تجربه ثابت می کند که انگیزه حفاظت، در میان تولید کنندگان محتوا کافی نیست و آن ها را برای حفاظت وابسته می سازد. در واقع، مرحله نخست حفاظت، اجبار در انجام آن توسط انواع مختلف سازماندهی، اجرای توسط اهداف متفاوت، انگیزه ها و حتی روش متفاوت است. وب، به عنوان زیر ساخت اطلاعاتی نمی تواند به طور اساسی مشکلات سازمانی را حل کند. از این رو، بایگانی کردن وب به عنوان فعالیتی مستقل از انتشار مورد نیاز است.

2-1-3- یک وظیفه غیر ممکن

سومین مقوله استدلال ها در مقابل بایگانی وب، از سوی افرادی مطرح می شود که نیاز به بایگانی وب را تصدیق می کنند، اما درباره امکان انجام آن شبهه دارند.

تردیده ها، یا در مورد اندازه وب است یا در موارد دیگر (تأکید بر خصوصی سازی، خاصیت، روشن فکر گرایی و موانع حق مؤلف) که بایگانی وب را به چالش می کشاند.

نخستین جنبه، متناسب به بیکرانی وب است که باید در رابطه با هزینه های ذخیره و ظرفیت ابزارهای خودکار برای گردآوری حجم زیاد اطلاعات مورد بررسی قرار گیرد. خطوط DSL فعلی و ظرفیت پردازش کامپیوترهای شخصی، خزش روزانه میلیون ها صفحه را امکان پذیر می کند. مقیاس میانگین بایگانی وب، در تناسب با مقیاس خود وب می باشد. حتی اگر تخمین دقیق آن مشکل باشد (داهن 1)، 2000؛ آگه 2000 (2)؛ دوبرا و فاینبرگ (3)، 2004) از منابع مختلف در می یابیم (4) که اندازه وب سطحی، به طور متداول در دامنه ده ها بیلیون صفحه است و این اطلاعات به شکل سایر نظام های اطلاعاتی پیچیده وب که - نمی تواند خزش کند (وب پنهان) - قابل دستیابی می باشد، البته به اندازه یا دو مرتبه بزرگ تر است.

بایگانی وب سطحی ثابت شده است طی یک دهه کامل توسط بایگانی اینترنتی، سازمانی کوچک با سرمایه گذاری خصوصی کوچک شدنی است (کاهل 5)، 2002، 1997).

دلیل این امر این است که این میزان مشابهی از محتوا، به وجود آورندگان، ارزش قابل توجهی را برای ایجاد، حفظ و نگهداری و دسترسی بالا پرداخت می کنند. ذخیره سازی فقط قسمت کمی از

ص: 13

Dahn -1

Egghe -2

Dobra and Fienberg -3

4- منابع به اندازه نمایه موتورهای کاوش مستند شده اند (یاهو ادعا می کند که 20 بیلیون صفحه را نمایه می کند، گوگل می گوید که بیشتر نمایه می کند (بتل، 2005) در یک نگاه کلی اندازه بایگانی اینترنت 10 بیلیون صفحه است)، مطالعات اخیر بر اساس روش شناختی های نمونه گیری است (گلی و سیگنورینی، 2005).

Kahle-5

هزینه های انتشار وب را امروز در بر می گیرد. بر عکس، بایگانی اینترنتی فقط برای ذخیره سازی، با استفاده از فشرده سازی برای مثال خزش گر توسط الکسا اهدا شده است) و دسترسی ها هزینه را پرداخت می کند و مورد دوم پرداخت برای هر واحد محتواست، که بسیار کوچک تر از چیزی است که سرور اصلی می پردازد. این نتایج در میزبان کردن یک کپی برداری کاملاً گسترده از وب در مؤسسه ای واحد (کوچک) به طور محسوس ممکن است.

جنبه دوم، نگرانی های خصوصی سازی، مالکیت معنوی و موانع حق مؤلف است فقط توجه داشته باشید که وب یک برنامه کاربردی انتشار غیر تجاری در اینترنت نیست. ارتباطات پنهانی برای رخ دادن در وب تصور نشده اند، اما درباره برنامه های کاربردی ارتباطات (مانند ایمیل و انتقال پیام) زمانی که این کار انجام می شود (لیوگ و فیشر (1)، 2003) همیشه احتمالات برای حفاظت از آن ها (که به طور وسیع می شود) به وسیله ورود به سیستم و اسم رمز وجود دارد. از این رو، فضاهای حفاظت شده به عنوان بخشی از وب عمومی مورد بررسی قرار نگرفته اند و بنابراین نباید در بایگانی های عمومی حفاظت شوند. این طرح طبیعی از جو خصوصی / عمومی در اینترنت به وسیله روشی که خزش گرها اجرا می کنند، تقویت می شود (به وسیله دنبال کردن پیوندها) به این معناست که صفحه ها و سایت ها به داشتن درجه معینی از پیوند درونی برای کشف شدن و تصرف شدن نیاز دارند. بقیه، اجزای غیر متصل وب هستند (برادر و همکارانش (2)، 2000) که به طور طبیعی از خزش گرها حذف می شوند. سایتی می تواند از این استفاده کرده و مرزهای بیشتری را برای شامل شدن در مجموعه (بیش از یک لینک درونی) برای محدود کردن تصرف بخش های مرئی تر تنظیم نماید.

در رابطه با وضعیت قانونی بایگانی وب، به وضوح، موقعیت های گوناگونی در هر کشور وجود دارد و این یک فضای در حال نمو می باشد. کشف این جنبه ها، فراتر از دامنه این کتاب است که در کتاب چارلزورث (2003) (3) به آن ها اشاره شده است. توجه داشته باشید که محتوای منتشر شده در وب غیر تجاری است چه توسط تبلیغات بر روی سایت ها یا بوسیله اشتراک پرداخت شود.

برای تمام موارد، بایگانی های وب، حتی با دسترسی پیوسته باید شرایط غیر رقابتی را با وبگاه های اصلی پیدا کنند و این کار می تواند در خصوص محدودیت های دسترسی به محتوا انجام شود (برای مثال همان طور که توسط تولیدکننده در متن فایل روبات ها گفته شده است). داشتن یک دوره ممنوعیت، قابلیت های کمتری را نشان می دهد (جست و جوی سایت و تعاملات پیچیده) و همچنین عملکردهای سطح پایین (سرعت دسترسی به محتوا). بنابراین، استفاده از بایگانی وب برای دسترسی به محتوا، زمانی انجام می شود که دسترسی اصلی امکان پذیر نباشد و درآمد بازدهی داشته باشد. در این صورت، برای ناشر اصلی استفاده از بایگانی وب تهدیدی محسوب نمی شود (این موضوع را در لایمن 2002 (4) ببینید). در مقابل، بایگانی وب می تواند به طور قابل توجهی برای به وجود آوردن سایت، حفظ بار محتوای منسوخ

ص: 14

Lueg and Fisher -1

.Broder et al -2

Charlesworth -3

Lyman -4

(قدیمی) را کم کند امکان تمرکز بر روی محتوای فعلی را می دهد. حتی در این موقعیت، نویسندگان و ناشران ممکن است درخواست کنند که مطالب شان از بایگانی های قابل دسترس عموم برداشته شود. درخواست همچنین می تواند توسط شخص سوم به دلایل مختلف انجام شود. چگونه بایگانی های عمومی وب به این درخواست ها پاسخ خواهند داد؟

توصیه هایی در این زمینه در ایالات متحده پیشنهاد شده است. جدول 1-1 را ببینید، (آبویس 2002) (1).

عکس

بایگانی شبکه وب: مباحث و روش ها 15

(قدیمی) را کم کند امکان تمرکز بر روی محتوای فعلی را می دهد. حتی در این موقعیت، نویسندگان و ناشران ممکن است درخواست کنند که مطالبشان از بایگانی های قابل دسترس عموم برداشته شود. درخواست، همچنین می تواند توسط شخص سوم به دلایل مختلف انجام شود. چگونه بایگانی های عمومی وب به این درخواست ها پاسخ خواهند داد؟
توصیه هایی در این زمینه در ایالات متحده پیشنهاد شده است. جدول 1-1 را ببینید، (آبویس 2002).

جدول 1.1

درخواست	توصیه
درخواست مدیر وب توسط یک وبگاه خصوصی (غیردولتی) به دلایل خصوصی سازی، افترا یا خجالت	<p>1- آرشیویست ها باید یک سایت به روش سلف سرویس ایجاد کنند و مالکان سایت می توانند مطالبشان را با استفاده از آن حذف کنند که براساس استفاده از معیار پروتکل رویوت متن می باشد (با این قابل متنی می توان، میزان دسترسی موتور جست و جوگر به محتوای یک سایت را کنترل کرد).</p> <p>2- درخواست کنندگان ممکن است بخواهند با دلیل مدرک مالکیت شان را با تغییر با افزودن یک پروتکل رویوت متن، در سایتشان اثبات کنند.</p> <p>3- این کار به بایگان ها اجازه می دهد تا مطمئن شوند که مطالب گذشته بیش از این جمع آوری نخواهند شد یا قابل دسترسی نیستند.</p> <p>4- این درخواست، عمومیت نخواهد داشت از این رو آرشیویست ها باید کپی های تمام درخواست های حذف شده را نگه دارند.</p>
درخواست های پاک شده شخص سوم براساس بیانیه حق مؤلف هزاره رقمی سال 1998 می باشد.	<p>1- آرشیویست ها باید در تلاش برای بررسی اعتبار شکایت ها با کنترل آنها باشند که آیا صفحه ها اصلی ثبت شده اند و اگر مناسب است بر روی درخواست نسبت به سایت اصلی از نظر قانونی اعمال حکم می شود.</p> <p>2- اگر شکایت معتبر باشد، آرشیویست ها باید موافقت نمایند.</p> <p>3- آرشیویست ها برای ایجاد درخواست های عمومی از بیانیه حق مؤلف هزاره رقمی، از طریق اثرات ناامیدکننده¹ و اخطار به جست و جوکنندگان زمانی که صفحه ها درخواست شده حذف شده اند، تلاش خواهند کرد.</p> <p>4- بایگان ها به مدیران سایت های مذکور از طریق ایمیل، اخطار خواهند داد.</p>
درخواست های پاک شده شخص سوم براساس شکایت های غیر از بیانیه حق مؤلف هزاره رقمی شخصی و (شامل علائم تجاری و رازهای تجاری در تولید).	<p>1- آرشیویست ها برای بررسی اعتبار شکایت با کنترل این که آیا صفحه ها اصلی ثبت شده اند، تلاش خواهند کرد، تا اگر مناسب است بر روی درخواست در رابطه با سایت اصلی از نظر قانونی اعمال حکم کنند.</p> <p>2- اگر صفحه ها اصلی پاک شده اند و آرشیویست ها تعیین کند که پاک کردن آنها از سرورهای عمومی، مناسب است، آرشیویست ها صفحه ها را از سرورهای عمومی شان بر می دارند.</p> <p>3- آرشیویست ها برای ایجاد این درخواست های عمومی به وسیله اثرات ناامیدکننده و اخطار به جست و جوکنندگان وقتی صفحه های درخواست شده برداشته شده باشد، تلاش می کنند.</p> <p>4- آرشیویست ها به مدیران سایت های مذکور از طریق ایمیل، اخطار خواهند داد.</p>

جدول 1.1

ص: 15

Ubois -1

درخواست‌های پاک شده شخص سوم براساس اعتراض به محتوای بحث انگیز باشد (مثل مباحث سیاسی، مذهبی و عقاید دیگر)	همان‌طور که در مجموعه قوانین کتابخانه بیل آ طبقت ذکر شده کتابخانه‌ها باید مطالب و اطلاعات را تهیه کنند و تمام نقطه نظرات درباره مباحث جاری و تاریخی ارائه دهند. مطالب نباید ممنوعیت (انتشار) داشته باشد یا به علت عدم تصویب یا عدم رضایت طرفداران پاک شود. از این رو، آرشیویست‌ها نباید به‌طور کلی به این درخواست‌ها عکس عمل نشان دهند.
درخواست‌های پاک شده شخص سوم براساس اعتراض به افشاء داده‌های شخصی که به‌طور محرمانه تهیه شده است.	گاهی اوقات، داده‌های محرمانه افشاء شده به‌وسیله یک طرف به طرف دیگر ممکن است در نهایت توسط یک شخص سوم عمومیت پیدا کند. برای مثال، اطلاعات پزشکی که به‌طور محرمانه تهیه شده، گاهی اوقات عمومیت پیدا می‌کند. وقتی عملکردهای شرکت‌های بیمه یا عملکردهای پزشکی متوقف می‌شود، این درخواست‌ها به‌طور کلی به معنای درخواست‌هایی تلقی می‌شوند که توسط نویسندگان یا ناشران داده‌های اصلی ایجاد می‌شود.
درخواست توسط دولت	آرشیویست‌ها بهترین تاثیر مرتبط با قبول قابلیت کاربرد قرار صادره از دادگاه، به‌کار می‌برند. بعلاوه، همان‌طور که طبق قوانین کتابخانه بیل در ذکر شد کتابخانه‌ها باید سانسور عقاید در اجرای تعهداتشان برای تهیه اطلاعات و آگاهی حقیقی را به چالش در آورند.
درخواست‌های دیگر و شکایت‌ها، مباحث حقوقی اصولی را در بر می‌گیرد.	درخواست‌های دیگر و شکایت‌ها، مباحث حقوقی اصولی را در بر می‌گیرد. اینها براساس مورد به مورد توسط بایگانی و مشاوران آنها انجام می‌شوند. کنترل و ایجاد مجدد وبگاه‌ها براساس تغییر مالکیت می‌باشد. این توصیه‌ها می‌توانند با محیط‌های قانونی دیگر وفق داده شوند تا جایی که استفاده مجدد از ساز و کارهای عملکردی مهم ایجاد گردند (ارتباطات از مالک سایت از طریق استفاده گسترده از استاندارد متن روبوتس و همترازی روی آن چه که بر سایت اصلی در ادعای سوم انجام شده است).

نیاز برای درک بهتر همزیستی بین به‌وجود آورندگان سایت و بایگانی‌های وب وجود دارد تا آنجا که بتوانند با عنایت به حقوق ایجاد کننده اطمینان پیدا کنند که از حافظه می‌تواند محافظت شود. اما این نیز بخشی از فرآیند تکامل رسانه وب است.

در مجموع، استدلال‌ها در برابر ضرورت و همچنین امکانات بایگانی وب است. جای تعجب نیست که، از نظر ما، در مغایرت با نقش مرکزی وب در خلق فرهنگ و انتشار آن و همچنین براساس طبیعت مطلقش، ایجاد شده است. فصل ۲ بینش بیشتری درباره چگونگی اهمیت بایگانی‌های وب برای پژوهش در بسیاری از حوزه‌ها را فراهم می‌کند.

در اینجا، سعی می‌کنیم نشان دهیم که در صورتی که چالش‌های جدی و مهمی را برای عملکردهای سنتی مطرح کردیم، بایگانی وب امکان‌پذیر می‌شود و یکی از موارد اصلی در برنامه حفاظت از میراث فرهنگی امروز است.

۳- ویژگی‌های وب برای حفاظت

وب خصوصیات مهمی دارد که هر تلاش حفاظتی باید درباره آن انجام شود. ما آنها را در این بخش در

نیاز برای درک بهتر همزیستی بین به‌وجود آورندگان سایت و بایگانی‌های وب وجود دارد تا آنجا که بتوانند با عنایت به حقوق ایجاد کننده اطمینان پیدا کنند که از حافظه می‌تواند محافظت شود. اما این نیز بخشی از فرآیند تکامل رسانه وب است.

در مجموع، استدلال‌ها در برابر ضرورت و همچنین امکانات بایگانی وب است. جای تعجب نیست که از نظر ما، در مغایرت با نقش مرکزی وب در خلق فرهنگ و انتشار آن و همچنین بر اساس طبیعت مطلقش، ایجاد شده است. فصل ۲ بینش بیشتری درباره چگونگی

اهمیت بایگانی های وب برای پژوهش در بسیاری از حوزه ها را فراهم می کند.

در اینجا، سعی می کنیم نشان دهیم که در صورتی که چالش های جدی و مهمی را برای عملکردهای سنتی مطرح کردیم، بایگانی وب امکان پذیر می شود و یکی از موارد اصلی در برنامه حفاظت از میراث فرهنگی امروز است.

3-ویژگی های وب برای حفاظت

وب خصوصیات مهمی دارد که هر تلاش حفاظتی باید درباره آن انجام شود. ما آن ها را در این بخش در

ص: 16

زوایای مختلف مورد بررسی قرار می دهیم. ابتدا، کاردینالیتی وب است یعنی اینکه چه تعداد نمونه از هر قسمت محتوا موجود است. و دوم اینکه وب به عنوان یک نظام انتشاراتی فعال و آخرین مورد وب به عنوان یک محصول فرهنگی جهانی، با طبیعت مافوق رسانه ای و طبیعت نشر آزاد بودن آن مورد بررسی قرار گرفته است.

3-1-1- کاردینالیتی وب

نخستین سؤال که برای حفاظت از محصولات فرهنگی عنوان شده است، کاردینالیتی بودن آن است. تعداد مواردی که هر اثر در حال توزیع شدن است. بایگانی ها و موزه ها معمولاً، با محصولات منحصر به فردی سر و کار دارند و حتی اگر در برخی موارد چندین قالب وجود دارد که کپی یا نشانه هایی از یک پیکر تراشی واحد، نقاشی یا اثر عکاسی است.

بر عکس، کتابخانه ها تقریباً موارد غیر منحصر به فرد را در مجموعه چاپی نگهداری می کنند (حفاظت از نسخه خطی، از این دیدگاه نزدیک به عملکرد بایگانی می باشد) منحصر به فرد بودن دارای اهمیت اجتماعی و نمادین عمیق است (بنجامین 1963) (1). همچنین، اثر بسیار زیاد آشکاری بر عملکردهای حفاظتی دارد. کتابخانه ها همیشه یک فرصت ثانویه برای یافتن کتاب های چاپ شده بعد از انتشارشان دارند.

چنین تخمین زده می شود بیشتر از 20 میلیون کتاب برای 30/000 ویرایش، بین 1455 و 1501 به چاپ رسیده است (فبوره و مارتین، 1976) که به این معنا که به طور میانگین نخستین دوره کاردینالیتی، متجاوز از 650 بوده است. کاردینالیتی مستلزم این است که حفاظت با تأخیر معینی بعد از انتشار رخ دهد، همان طور که کپی های متعدد برای یک دوره زمانی حتی در غیاب حفاظت فعال باقی می مانند. همچنین، یک سطح طبیعی از افزونگی یک ویژگی در یک نظام را موجب می شود که کتابخانه متفقاً انجام می دهند. با استفاده از داده های یکی از بزرگ ترین پایگاه داده های کتابشناختی (worldcat) لایوه و شانفلد (2005) (2) سه ردیف توزیع کاردینالیتی اثر منتشر شده در کتابخانه ها را کشف نمودند که از Worldcat استفاده کرده است (تقریباً 20/000 در آمریکای شمالی): 37 درصد فقط یکبار نگهداشته شده اند، 30 درصد، 2-0 بار و 33 درصد بیشتر از 5 بار نگاه داشته شده اند.

زمان و افزونگی (تکرار اطلاعات میان فایل های گوناگون) دو مزیت قابل توجه از یک چشم انداز حفاظتی هستند که یکدیگر را تقویت می کنند. آن ها همیشه وجود ندارند. تولید مجدد نسخه های خطی در رابطه نقایص آن برای قرن ها قبل از اختراع چاپ اقدام شده بود، از این رو، اکنون حتی زمانی که چندین نسخه (واقعاً تعدادی) گوناگون وجود دارد. کتابداران بزرگ ترین کتابخانه قدیمی اسکندریه (3) را تشکیل می دهند که از کپی برداری نسخه های خطی استفاده می کردند که به دورن شهرها انتقال یافته بودند اما آن ها

ص: 17

Benjamin -1

Lavoie and Schonfeld -2

Alexandria -3

نسخه اصلی را نگهداری می کردند (کن فورا، 1989).

ترجمه و جمع آوری، تفسیر و توضیحات و حاشیه نویسی، غالب اوقات بنیاد و پایه اصلی برای تولید مجدد متن به جای حفاظت قابل اعتماد بوده اند که برای زیان و ضررهای اجتناب ناپذیر اضافه شده بود که مستلزم کپی برداری دستی بودند. بیشتر کپی برداری های اصولی از متون، اغلب به دلایل خارجی ایجاد می شدند مانند وقتی متون یونانی، اساساً در موقعیت اختراع یک نوشته جدید، (حروف کوچک) در دوران امپراطوری رم در قرن نهم کپی برداری شده اند، تثبیت و انتقال آن ها به شکلی که ما امروزه می شناسیم خواهد بود.

کپی برداری آینده به طور قابل توجهی، شرایط را در این رابطه تغییر می داد که آن محتوا را به حالت تثبیت کرد زمانی که توزیع گسترده تر آن را مجاز کرد (آیزنشتاین 1979 (1)؛ فبوره و مارتین (1976) (2). همچنین، با افزایش قابل توجهی از کاردینالیته آثار، راندمان حفاظت را بدون سابقه کرد. در جایی بر آورد شده است که یکی از 40 اثر شناخته شده از دوران قدیم، حفظ شده است (و کمتر اگر آثار ناشناخته را در نظر بگیریم). راندمان حفاظت به بیشتر از یکی از دو تا در قرن هفدهم در فرانسه و نزدیک 80 درصد یک قرن بعد از آن (استیوالز 1965) (3) بالا رفته است و برای یک مؤسسه واحد، کتابخانه سلطنتی نیز بعد از تقویت سپرده گذاری قانونی توسط فرانسیس ل (4) در سال 1537 (استیوالز، 1961؛ بالیه، 1988) (5).

امروزه، حفاظت از کارهای چاپ شده، در بیشتر کشورها به کارآیی و رشد مؤثری دست یافته است؛ از نقطه نظر عملی و سازمانی که با ثبات مطالب چاپ شده و همچنین کار دینالیته، مجاز شده اند.

هر آن چه که بود، کاردینالیته در محصولات فرهنگی حداقل از ایجاد تا دسترسی یکپارچه می باشد. این تنها مورد در وب نیست. کاردینالیته محتوایی وب ساده نیست، بلکه هر یک (چند جزئی) است. همان طور که منبع محتوا معمولاً یک سرور منحصر به فرد است شخص می تواند به طور محسوس نماید که کاردینالیته اش مانند آثار هنری و نسخه های خطی، یکی می باشد. در واقع، همان آسیب پذیری را نشان می دهد، حتی توسط این حقیقت که محتوا به تولید کننده خود وابسته است، افزایش می یابد. از طرفی دیگر، دسترسی و همچنین کپی های محتوای وب می تواند از نظر مجازی، نامحدود باشد. این اختلاف میان دو کاردینالیته های وب، ما را به سمت این مفهوم مهم از منبع وب هدایت می کند. هر منبع، یک منبع منحصر به فرد (سرور وب) و یک شناسه منحصر به فرد دارد، اما می تواند از نظر مجازی به طور نامحدود و با درجه های گوناگون برای برنامه ریزی های گوناگون تولید شود.

از دیدگاه حفاظتی، هر منبع دو خصوصیت مهم دارد:

نخست اینکه به طور دائمی به منشأ انحصاری اش برای موجودیت وابسته است. این کار یک تفاوت قابل توجهی با انتشار ایجاد می کند، جایی که مدیران انتشار، فقط یک بار به آن نیاز پیدا می کنند و بعد

ص: 18

Eisenstein -1

Febvre and Martin -2

Estivals -3

François 1er -4

از آن، کتاب‌ها به وجود می‌آیند. دوم اینکه سرورهای وب می‌توانند محتوا را برای هر نوع منبع، مناسب سازند و آن را در هر زمانی برای یو.آر.ال مشابه متفاوت می‌سازد. وب از این دیدگاه، یک ظرف محتوی فایل‌های ثابت نیست، اما یک جعبه سیاه با منابعی است که کاربران فقط نمونه‌هایی را به دست می‌آورند.

همان‌طور که کریشنامرتی (1) و رکسفورد (2) درباره پروتکل وب توضیح داده‌اند:

یک روش برای ادراک پروتکل، این تصور است که منبع سرور حاوی جعبه‌های سیاه با نمایش منابعی باشد که توسط یو.آر.ال‌ها معنا شده‌اند. منبع سرور اصلی، شیوه درخواست برای منبع مشخص شده را به وسیله یو.آر.ال درخواست می‌کند. دریافت مشترک از خواندن یک منبع از یک فایل و نوشتن پاسخ برگشت به سرویس‌گیرنده، دور از دید جعبه سیاه، مجزا و مختصر شده است. این نگرش، مفهوم یک منبع را عمومیت می‌بخشد و آن را از پاسخ ارسال شده به سرویس‌گیرنده تفکیک می‌نماید. درخواست‌های مختلف برای یو.آر.ال‌های مشابه، می‌تواند ناشی از پاسخ‌های متفاوت باشد و به عوامل مختلفی بستگی دارد. فیلدهای بالایی درخواست، زمان درخواست با تغییرات برای منابعی که ممکن است رخ دهند (کریشنامرتی و رکسفورد، 2001) حفاظت وب، منابع را طبق کاردینالیتی دوگانه (به ظاهر مهم و در واقع درست) مورد بررسی قرار می‌دهد و این کار مستلزم چندین استنباط است. نخست اینکه چون از نظر مجازی تعداد نامحدود کپی برداری می‌تواند به آسانی ایجاد شود، شخص می‌تواند ادارک گمراه‌کننده‌ای داشته باشد که بایگانی فعال وب برای حفاظت مورد نیاز نیست. از این رو، تعدد نمونه‌ها، به طور گسترده مخفی‌اند و به یک منبع تکی بستگی دارند. هر زمانی که لازم باشد. (سرور) کد می‌تواند برچیده و روزآمد شود از این رو برای یک بایگانی فعال مورد نیاز می‌باشد.

استنباط دوم این است که بایگانی‌های وب می‌تواند فقط برخی موارد در منابع را به طور بالقوه درجات گوناگونی از میان آن‌ها را به تصرف در آورد (3). این مورد زمانی رخ می‌دهد که محتوا برای یک مرورگر خاص یک زمان معین یا یک موقعیت جغرافیایی معین یا زمانی که محتوا با هر کار بر وفق داده شده است، مناسب گردد. همچنین، در بخش بعد خواهیم دید که وب در حقیقت یک سیستم نشر فعال است و از این رو، تفاوت پاسخ‌ها در واقع جنبه‌ای مهم برای بررسی است، به خصوص وقتی که بایگانی انجام می‌شود.

3-2-وب به عنوان سیستم نشر فعال

وب، نوعی برنامه کاربردی نشر اصلی در اینترنت است همچنین، به طور اساسی شامل ترکیب سه

ص: 19

Krishnamurthy -1

Rexford -2

3- توسعه پویای صفحه‌ها برای ایجاد یگانگی در طراحی و معماری نیز در کل سایت استفاده شده است (دستگاه‌های شناوری و مانند آن). استفاده از تمپلیت‌ها به طور یکسان نگاه کردن به صفحه‌ها را آسان کرده و تغییر طراحی توسط تمپلیت‌ها را آسان‌تر از صفحه‌های انفرادی می‌کنند برآورد شده است که تمپلیت‌های مبتنی بر صفحه‌ها 40 درصد تا 50 درصد از صفحه‌ها را نمایش می‌دهند (جیبسون و دیگران 2005).

استاندارد می باشد: 1) URL (لی - برنرز، 1994) که فضای نام گذاری شده را برای یک شیء تعریف می کند (1)؛ 2) HTTP (فیلدینگ و همکارانش (2)، 1994) پروتکل تعامل سرویس دهنده - سرویس گیرنده را با استفاده از فرآیندها در هسته اش تعریف می کند؛ و 3) HTML (برنرز-لی و کونولی (3)، 1995)، نوعی (4) SGML DTD (تعریف نوع داده) که ارائه صفحه آرایی در مرورگرها را معین می نماید. اجرای این سه استاندارد، هر کامپیوتر متصل به اینترنت را برای وارد شدن به سیستم نشر، قادر می سازد. شبکه سرویس دهندگان وب، نوعی سیستم اطلاعاتی منحصر به فرد را تشکیل می دهد که می تواند در هر حالتی برای تولید، روزآمد شدن، و نشر محتوا در حالتی که کامپیوترهای جدید اجازه می دهند، به کار رود.

در مقایسه با وسایل انتشاراتی دیگر، انقلاب در نشر، گسترش امکانات در تمام جهات ممکن برای تولید، سازماندهی، دستیابی، و ارائه محتوا را نشان می دهد. برای مثال، پیوندها را مورد بررسی قرار دهید:

شخص می تواند استدلال کند این فقط شکل جدیدی از یک ارجاع است که پیش از زمانی که برای نخستین بار نوشته شده به وجود آمده است (5). اما حقیقت این است که روشی که به وسیله قطعه قطعه کردن محتوا به تکه های نشانی پذیر کوچک تر و توجه کلی به برنامه ریزی خاصیت انتقال از طریق شناوری دسترسی به محتوایی که تغییرات عمیقی را در نوشتن و همچنین در خواندن پیدا کرده است، به دلیل تغییرات وب در روش ارجاع قابل تعقیب قانونی است (آرسث (6)، 1997؛ لندو (7)، 1997؛ بولتر (8)، 2001).

این حقیقت، که محتوا تنها بر روی سیستم و با دقت بیشتری بر روی سرورهای ناشران موجود است، به انتشار دائم ایجاد کننده بستگی دارد. یک کتاب می تواند بعد از ترک چاپخانه به طور مستقل از ناشر باقی بماند، اما محتوای وب هیچ موجودیتی فراتر از سرور اصلی اش نخواهد داشت (به استثنای ساز و کارهای حافظه با سرعت بالای ناپایدار (9) هافمن و بیومونت، 2005). انتشار دائم، کنترل واضح را گسترش می دهد که ایجادکنندگان بر روی محتوا دارند. آن ها می توانند، با وب، در هر زمانی تغییر کنند، به روز شوند، و در زمان واقعی مواردی را از انتشار پاک کنند. علاوه بر این، تولیدکنندگان وب، از نظام اطلاع رسانی وب (10) استفاده می کنند که بتواند اطلاعات را از هر نوع نظام اطلاعاتی موجود (پایگاه داده ها، مخزن اسناد، برنامه های کاربردی، و مانند آن) ترکیب، مجتمع و سازماندهی مجدد کند. از این رو، وب یک فضای اطلاعاتی ثابت نیست، بلکه یک فضای نشر فعال است که ناشی از اثرات یک مجموعه آمیخته شده از نظام های اطلاعاتی فعال می باشد.

ص: 20

1- این استاندارد مهم ترین در میان سه مبتکر وب است (لی - برنرز و فیشتی، 2000؛ گیلز کیلیو، 2000؛ چنان چه وب را در موقعیت دسترسی جهانی قرار داده است که کلاً منبع سندی قابل دسترس در اینترنت است.

Fielding -2

Connolly -3

SGML DTD -4

5- برای مقایسه استنادهای علمی سنتی و این که چگونه می تواند برای ارزیابی علمی استفاده شود اینگورسن (1998) را ببینید بچورن بورن و اینگورسن (2001) تحلیل انتقادی آن در توال (2001)، توال و هریس (2004) و توال (2006).

Aarseth -6

Landow -7

Bolter -8

Hofmann and Beaumont -9

(Web information systems (WIS -10

از این رو، بایگانی وب نخست به جدا کردن محتوا از نشر ثابت ایجادکنندگان اصلی اش نیاز دارد، و دوم اینکه باید مطمئن شود که محتوا می تواند از عدم پذیرش و تکامل جاری وب، عدول کند.

قبلاً به کپی برداری و بایگانی محتوا در یک زیر ساخت مجزا نیاز بود (مطالب زیر، (1) فصل 3 و روشه، 2006 را ببینید). مورد آخر مستلزم حفاظت فعال از محتوای وب (فصل 8، دی (2) 2006 را ببینید) برای رفع وابستگی از اجزای سیستم های گوناگون (پروتکل ها، به فرمت های دیجیتال، برنامه های کاربردی، و نظیر آن) و اجتناب از منسوخ بودن اصول فنی آن هاست. حفاظت از وب، این نیازها را در کل برای حفاظت از اصول فنی فعال با اشیای رقومی به اشتراک گذاشته است، اما جداسازی از ایجاد کننده نشر دائم، برای حفاظت از وب مشخص است.

اما رفع هر گونه وابستگی از سرور اصلی، مستلزم این است که از قابلیت های گوناگون و شیوه تعاملی وب، بایگانی وب بتواند فقط تعداد کمی را حفظ کند. هزینه هایی برای جداسازی از شبکه اصلی نظام های اطلاعاتی وب وجود دارد.

قابلیت هایی که بر بخش سرویس گیرنده ها اجرا می شوند، آن هایی هستند که شخص می تواند به طور معقولانه به حفظ آن امیدوار باشد. دامنه قابلیت هایی در کد صفحه و کد فایل مربوط جاسازی شده اند که به میل سرویس گیرنده اجرا می شوند و بیشتر اوقات بر روی نسخه های بایگانی قابل اجرا هستند؛ اما این قابلیت ها که توسط کد و یا با اطلاعات سرور تهیه شده اند جاسازی نمی شوند. جنبه های سندی مطالب اصلی هستند که گم شده اند (مانند انواع تعاملات که شخص می تواند بر روی یک ویدئو ثبت نماید)، اما این کار فقط می تواند برای تعداد محدودی از صفحه ها و نقطه نظرات معین و شرایط خاص انجام شود (کریستسن - دالسگاد (3) ، 2001 ؛ بروگر، 2005) (4).

3-3-وب به عنوان یک محصول فرهنگی

علاوه بر یک نظام نشر فعال، وب یک فضای اطلاعاتی با مشخصات خاص است. واژه وب در این مضمون، یک محصول فرهنگی رقومی گسترده را برگزیده است (لایمن و کاهله، 1998) که می تواند با حقایق زیر مشخص شود:

ص: 21

Roche -1

Day -2

Christensen-Dalsgaard -3

4- نقطه نظر طراح وبگاه در این مورد نیز جالب است. در دابری و دیگران (2002). جلیس هودج پیشنهاد می کند، سایت هایی که او در حال طراحی است بایگانی شوند: - درخواست برای پروپوزال؛ - بیان هدف و دلیل استفاده؛ - توصیف استفاده از مضمون (مثال های مورد نیاز)؛ - توصیف استفاده کنندگان واقعی و مورد هدف؛ - جایگزین های ثابتی که به اندازه کافی دیدگاه و احساس را احاطه می کنند؛ - مثال هایی از چند راه مهم در سایت؛ - توصیف فناوری هایی که به کار برده شده و یا حمایت می کنند؛ و - هر ماژول مرتبط مثل پویا نمایی های فلش، فیلم ها، PDF ها، و مانند آن.

- از هر محلی که متصل به اینترنت است قابل انتشار و دسترسی (معمولاً مجانی) می باشد؛

- به عنوان یک فوق رسانه با استفاده از پیوندهای مستقیم و قابل تعقیب قانونی بین قطعات محتوا ساخت

یافته است (1)؛

- نه تنها شامل متن است، بلکه شامل ترکیبی از تصاویر، صداها، و محتوای متنی نیز می باشد؛ و

- ناشی از یک تألیف و تصنیف باز و توزیعی می باشد (2).

اگر چه وب، این کارها را به طور بسیار گسترده انجام می دهد، از اشکال قبلی انتشار هم می نماید (3) (کراستون و ویلیامز 1997) (4)؛ اریکسون و آیهلستروم، 2000؛ شفرد و پولانی، 2000). همچنین، مواردی جدید را ابداع می کند. برای مثال، بلاگ ها با تلفیق ساده گسترده ای منتشر می شوند (حتی مهارت های فنی که برای سایت های عادی لازم است و بیش از اینها مورد نیاز نیستند). و یک مدیریت مرجع قدرتمند (شامل مراجع معکوس یا آگاه سازی از استنادها با استفاده از برگشت پینگ) و سهولت در به روزرسانی، افزودن توضیحات و حذف محتوا، همه اینها ناشی از یک نشر آزاد و توضیحات شخصی به وسیله ده ها میلیون نفر بوده است (5).

این خصوصیات وب، به عنوان یک فوق رسانه توزیعی، به طور آشکار و ثابت در یک مقیاس سراسری، تألیف و فراهم شده است که بایگانی وب بتواند فقط، به حفاظت از جنبه های محدود شده محصولات فرهنگی موجود و بزرگ تر دست یابد.

اتصالات درونی محتوا، یک کیفیت مهم در وب است که زمانی که بایگانی می گردد، مباحثی از آن ها ایجاد می شود، اما به عنوان یک نتیجه عمومی، نشان داده است که بایگانی همیشه، بهترین نوع گزینش (به گزینی) را در برنمی گیرد، حتی اگر در موارد انتخاب دستی و سایت به سایت باشد. از این مباحث برای بایگانی های بزرگ و وسیع اجتناب می شود، تا جایی که ممکن است، قطع تداوم اطلاعاتی که وب ارائه می دهد (لایمن و دیگران 1998). یا برای تعریف یک هدف تحلیلی خاص برای زمینه یابی تصمیم های گزینش انجام می شود (بروگر، 2005).

اما در عمل، اجرای خزش، در ارتباط با موارد اولویت دار و سیاست گزینش دستی، اجزای بایگانی وب را در بر می گیرد که همیشه فقط یک برش در فضا و زمان وب اصلی خواهند بود.

چگونه این نمونه های معنی دار و جایگزین را در وب بزرگ تر ایجاد می کنید؟ چه استنباطی هایی

ص: 22

1- آبرون و مک کرلی (2003) نشان می دهند که یک سوم پیوندهای مستخرج از بیلیون ها صفحه از نقطه جهت مشابه می خزند، یک سوم از عرض، بالا یا پایین در سلسله مراتب جهت ها از همان سایت و یک سوم از سایت خارجی.

2- نویسندگی محدود به تعدادی از افراد نیست بلکه از طریق ده ها و صدها فرد توزیع شده است. مثلاً تخمین زده شده است که در مورد فرانسه، گره های نشر، شخص یا ساختاری که منتشر می کند (ویراستاران نه نویسندگان) با انتشار در وب سه برابر اندازه توسعه یافته است: در حدود 5000 ناشر یا ساختار دهنده نشر به پنج میلیون سایت و سایت شخصی (منبع: انجمن فرانسوی دسترسی به اینترنت). این موارد

شامل وب نوشت ها نیست.

3- این در مورد چاپ نیز بود که برای مدت طولانی از دست نوشته ها و سازماندهی صفحه قبل از اختراع چاپ تبعیت می کرد (فبوره و مارتین، 1976).

Crowston and Williams -4

5- در مورد حفاظت از بلاگ ها انتلیج (2004) را ببینید.

در درک از آینده کاری که وب انجام خواهد داد دارید؟ تمام این سؤال‌ها باید زمانی که بایگانی و را به کار گرفته اید مورد بررسی قرار بگیرند. حتی تعریف کاری که «وب اصلی» انجام می‌دهد مجموعه تجربیات استفاده‌کنندگان از وب یا مجموع برنامه‌ریزی‌های محتوای وب است که باید از قبل در نظر بگیریم که بر مجموعه‌ای از محتواهای ثابت بیشتر در نظر گرفته می‌شوند.

خصوصیات دیگری که برای اندیشه و سازماندهی مجدد عملکردهای حفاظت سنتی باید انجام طبیعت نویسنده‌گی باز وب است. در حقیقت، این موضوع فیلتر کردن و حفاظت ساختار را بسیار مشکل می‌سازد که بر اساس ناشران و نویسندگان است. آن‌ها به اندازه بسیار قابل توجهی بر روی وب می‌باشند و برای شناسایی و ثبت مشکل‌اند. گاهی اوقات، اطلاعات تألیف و تصنیف، بر روی سایت موجود است، گاهی اوقات هم نیست، و گاهی اوقات در روش قابل اعتمادی نیز قابل دسترسی نیست. تنها اطلاعات ثبت شده (در یک روش کاملاً کنترل نشده و آزاد)، اطلاعات درباره کسانی که نام دامنه را برای مدیری DNS مجاز می‌کنند. اگر چه این اطلاعات مقادیر بزرگی برای کامل کردن مطالب بایگانی شده در وب هستند، و به طور قطع، برای تفسیر و استفاده مستقیم آسان نیست.

به عنوان یک محصول فرهنگی، وب، سبک متفاوتی از سازماندهی اطلاعات و الگوهای ساختاری متفاوت را برای استفاده در حفاظت ارائه می‌دهد. نشانه‌های پیوند محتواها و شناوری کاربران از ساختارهای طبیعی که بایگانی‌ها باید بیشتر وقت‌ها برای سازماندهی موارد نمونه‌های جمع‌آوری شده استفاده کنند. از این رو، خصوصیات وب به تغییر شکل و روش‌های حفاظتی عمیق نیاز دارد. رویکرد کل‌نگری برای بایگانی وب، آمادگی بیشتر برای انطباق با خصوصیات وب است، اما هر نوع بایگانی وب باید آن‌ها را در هسته روش‌های مشارکت دهد.

4- روش‌های جدید برای رسانه جدید

کتابخانه‌ها آرشیوها (بایگانی)، و موزه‌ها، روش‌های بسیار کارآمدی را با موضوعات مورد علاقه‌شان تطبیق داده‌اند که نقش مهمی در ساخت حافظه اجتماعی ایفا می‌کنند. اگر چه باید بیشتر فراگرفته شوند و بتواند برای حفاظت از وب مجدداً استفاده‌گردند. ماهیت وب و کیفیت‌های مورد نیاز، همان‌طور که دیده شد برای بررسی مجدد و تطبیق عملکردهای حفاظتی به ارث برده شده از این سنت طولانی در حفاظت از محصولات فرهنگی، فیزیکی می‌باشد. این بخش، یک بررسی عمومی از روش‌های جدید و رویکردهایی را نشان خواهد داد که باید برای حفاظت از وب مورد استفاده قرار گیرند (فصل 3 تا 8 مباحث را با جزئیات تهیه می‌کنند).

قبل از شروع مبحث روش شناختی، باید درباره استقرار بایگانی وب در زیرساخت های اطلاعات (1) به طور کلی، و به ویژه در اینترنت پرسش هایی مطرح شود.

بورگمن (2000)، تعریف کتابخانه رقومی جهانی را مطرح کرد (pp47sq) و تفاوت میان نگرش تکاملی و انقلابی در فناوری اطلاعات را توضیح داد:

«نگرش انقلابی، کتابخانه هایی است که رقومی هستند و با پایگاه های داده ها به وسیله شبکه های کامپیوتری پیوند یافته اند و در کل، می توانند آرایه ای از خدمات را فراهم نمایند که کتابخانه ها را از ریشه برخواهند کند. نگرش تکاملی این است که کتابخانه های رقومی سازمان هایی هستند که به تهیه محتواها و خدمات در شکل های گوناگون و فقط به عنوان پیش نیازهای مؤسسه ها ادامه خواهند یافت و سرانجام کتابخانه های مکملی خواهند بود که امروزه وجود دارند» (همان، ص 48).

او، تعریف میان گذر از تکامل - انقلابی را پیشنهاد کرده است: که «کتابخانه های رقومی یک گسترش، افزایش و یکپارچه سازی در نظام های بازبایی اطلاعات و مؤسسه های با اطلاعات چندگانه و کتابخانه ها که فقط یکی باشد را بیان می کند. محدوده قابلیت های کتابخانه های رقومی نه تنها شامل بازبایی اطلاعات است، بلکه ایجاد و استفاده از این اطلاعات می باشد».

موقعیت برای بایگانی های وب در جهتی که آن ها برای موجود شدن مورد بررسی قرار می گیرند و قبلاً فضای اطلاعات را ساخته اند، متفاوت است. همچنین، به طور آشکار قابل دستیابی اند. از روی سنجش و اندیشه، در این فضا، فقط به دروازه بان ها نیاز نداریم. همان طور که محدودیت های دسترسی فیزیکی وجود ندارند. در این زمینه، نقش بایگانی وب، در سازماندهی اطلاعات زیاد است.

کتابخانه های فیزیکی باید در هر دو سازماندهی فیزیکی و فکری اشیا ایجاد شوند و این طیف وسیعی از امکانات و انتخاب ها را میسر می سازد. همچنین، در حالی که دسترسی فیزیکی به محتوا را مدیریت می کنند، نقش میانجی گری اجتناب ناپذیری دارند. کتابخانه های رقومی، این نقش میانجی گری را به وسیله ایجاد محیط های همکاری و دانش وابسته به متن تحت تابع بنیانی جست و جو و دستیابی را گسترش می دهد (لاگوز و دیگران (2)، 2005).

بایگانی های وب در قسمت های مربوط به خود محتوایی بارگذاری شده با روابط تعبیه شده و قابل

ص: 24

1- مفهوم زیر ساخت به طور کل در استار و روهتلر (1994) با ابعاد مختلف تعریف شده است: - جاسازی (در هم جا دادن)، - شفافیت؛ - به دست آوردن یا دامنه (ساختار به یک حادثه عواحد یا عمل یک جانبه رسیده است)؛ - به عنوان بخشی از عضویت فراگرفته شده است (اعضای جدید برای عضو شدن نیاز به یک آشنایی دارند)؛ - پیوندهایی با کنوانسیون ها؛ - بر پایه پیاده سازی ساخته شوند؛ و - در حالت تفکیک قابل رؤیت باشند. این مفهوم در مضمون ساختارهای اطلاعاتی در بورگمن، (2000، 2003) بحث شده است

2-Lagoze et al.

تعقیب قانونی و ساختارهای اطلاعاتی غنی ایجاد شده که توسط میلیون ها نفر در سراسر جهان ویرایش می شود وقتی بایگانی های سنتی و کتابخانه ها نگرش سازمانی، شخصی و ابزارها را در این محتوا ایجاد می کنند (مدخل های شخصی و وب سازان و مانند آن)، فقط در این ویرایش جهانی وب شرکت می کنند. این کیفیت را به عنوان متخصصان حوزه تقلیل نخواهد داد اما آن ها را در یک تلاش سازمانی بزرگ تر قرار خواهد داد.

به عنوان بایگانی، وب آن ها مسئولیت های بیشتری دارند چون محتوا و مفهوم را تصرف کرده و تحت سیطره خواهند داشت و می توانند آزمایش هایی را در بهبود نقش منحصر به فرد قدیمی خود در سازمان دهندگان اطلاعات داشته باشند همچنین می توانند فقط برای تثبیت و محافظت از نمونه های شخصی از یک محصول فرهنگی جاودانی بزرگ تر به دست آیند.

این کار می تواند قانونی باشد زمانی که طبق خط مشی گزینش مناسب با نیاز جامعه کاربران برپا یا توسط اهداف پژوهشی معین شده اداره شود لچر (1) ، 2006؛ میزانس، 2006 ب). اما هزینه ها و محدودیت ها و همچنین امکانات فنی برای بایگانی هر دو در یک مقیاس بزرگ تر و در یک مسیر بی طرف نیز نیازمند بررسی روش های جایگزین در بایگانی وب می باشند این جایگزین در نقش خود معتدل تر ولی در دامنه جاه طلب تر است. نقش سازمان دهندگان اطلاعات برای به تصرف در آوردن محدود است و برای ساختار اصلی ایجاد شده توسط ویرایش میلیون ها نفر در سراسر جهان درست است.

در مشکل رسیدن به جامعیت، همچنین در بخش قبلی دیدیم که شخص می تواند حداقل تلاش در بی طرفی را برای تصرف محتوا با پیروی از ماهیت جمع آوری و توزیعی وب برای راهنمایی در تصرف و گسترش آن، تا جایی که ممکن است داشته باشد از این رو تلاش بر روی کمیت و موضوع مقیاس گذاری بر روی منابع فنی می باشد. این رویکرد، چندین کتابخانه ملی برای حوزه ملی و بایگانی اینترنتی در یک مقیاس جهانی داشته است.

هیچ یک از این ابتکارهای عملیاتی نمی توانند از نظر عمق و کیفیت محتوای بایگانی شده به تنهایی توسعه یابد. تلاش های گوناگونی به عنوان بخشی از یک بایگانی جهانی مورد بررسی قرار خواهد گرفت البته وقتی که اتصالات درونی میان بایگانی وب به عنوان اتصالات درونی بین سرورهای نشر از طریق وب سازماندهی شوند.

تنها با این کار، کاربران قادر به نفوذ به تمام این تلاش ها خواهند بود و به بهترین حافظه ممکن وب منتج خواهد شد. در این جهت، هر چه شراکت بیشتر مؤسسه ها و افراد مختلف وجود داشته باشد، بهتر می توانند مکمل یکدیگر باشند و زوایا، عمق و کیفیت های مختلف بایگانی های متفاوتی را ارائه دهند. اما این کار نیازمند این است که آن ها در برخی نقاط از یک شبکه بایگانی وب بزرگ تر، همکاری کنند. چنین شبکه ای باید بایگانی وب را پیوند دهد، به طوری که با یکدیگر نوعی فضای شناوری جهانی مانند خود وب را شکل دهند این امر فقط در صورتی ممکن است که آن ها در یک مسیر نزدیک به وب اصلی ساخته

ص: 25

شوند و آزادانه قابل دسترسی باشند کنسرسیوم (شرکت) حفاظت از اینترنت بین المللی (1) بر روی ایجاد و تنظیم زمینه هایی برای ایجاد کننده به وسیله توسعه استانداردها و ابزارهایی که ساخت این نوع بایگانی را تسهیل می کند عملکردهایی داشته است (برخی از آن ها در پایان این بخش توصیف خواهند شد).

دسترسی باز در خصوص مقررات و خط مشی است و در این برهه از زمان به عنوان یک بحث آزاد باقی می ماند.

بایگانی وب، به طور اختصاصی یا به عنوان کل می تواند در ایجاد زیر ساخت های اینترنتی مناسب باشد. آن ها از پروتکل ها و استانداردهای مشابه برای سازماندهی اطلاعات و ایجاد دسترسی به آن استفاده کنند. وب می تواند به طور طبیعی آن ها را در برگیرد چون آن ها کاملاً با آن سازگار هستند (2). از نقطه نظر زیربنایی بایگانی وب می تواند به آسانی موقعیتی را به عنوان مکمل ایجاد زیر ساخت های اینترنتی پیدا نماید. آن ها در حال فراهم کردن حافظه وب هستند که خود بخشی از وب است و اثر شدید منفی ماهیت ضروری ناپایدار انتشار وب را محدود می نماید.

شخص می تواند با چشم پوشی از این نقش ناراضی باشد شرایط این کار در نظر نگرفتن ارزش طبیعت توزیع گراو جامع بودن این رسانه است که آن را توجیه می کند.

4-2-4- فراهم آوری

اصطلاح فراهم آوری برای معانی فنی گوناگونی به کار می رود، مانند رسیدن محتوا به درون بایگانی. این اصطلاح شامل تصرف پیوسته و غیر پیوسته تحویل محتوا می شود که فرآیند انتخاب را پوشش می دهد و نه فرآیند درج و توسعه فردا را.

از دیدگاه فنی این مرحله از تعامل با تولید کنندگان و در صورت سنتی برای مؤسسه میراث حافظه، هر چیزی می تواند با شرایط در بایگانی وب جزئی می باشد. زیرا هیچ روش واحدی برای فنون انتشار گسترده وب کافی نیست گسترده سازی دامنه سازندگان و افزایش اندازه محتوا به درجه ای معین با خودکار سازی که در محیط وب ممکن می شود متعادل می شود. از این رو، موانع اصلی اکتساب ابزارهایی است که باید بر آن غلبه یابند و عدم توانایی پروتکل HTTP در ایجاد کپی دسته ای از محتوای سرور می باشد. سرورهای HTTP فقط می توانند تا زمانی که URI درخواست نماید، فایل به فایل تحویل دهند. این کار موجب کشف گذرگاه فردی برای هر فایل را در یکی از مباحث کلیدی در بایگانی وب

می شود.

در این بخش سه نوع روش فراهم آوری را بررسی می کنیم. چرا سه روش؟ چون فرآیند جمع آوری می تواند یا به عنوان یک سرویس گیرنده دور افتاده، در نزدیکی به خروجی سرور انجام شود یا به وسیله دسترسی مستقیم به فایل های سرور صورت گیرد (تصویر 1). گزینه نخست با خزنده بایگانی یا ماشین

ص: 26

وب باید از بایگانی کردن سایر آرشیوهای وب اجتناب ورزند و خود را با وب زنده (موجود) محدود کنند.

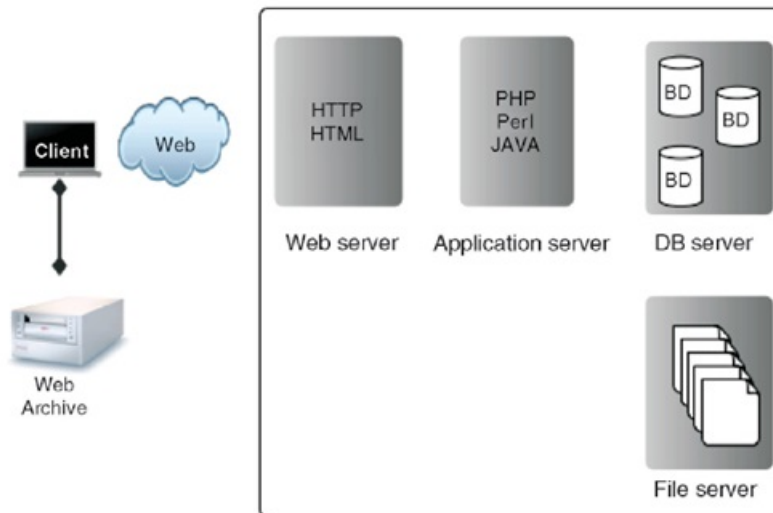
کپی کننده وبگاه انجام می شود، مشتق شده و در فناوری موتور جست و جوی سازگار و استنتاج می شود یک ابزار قدرتمند برای تصرف موقعیت سرویس گیرنده را فراهم می کند. روشه (2006)، توصیفی با جزئیات این ابزارها و کاربردها برای بایگانی وب ارائه می دهد. در این مقاله فقط یک بررسی عمومی از این فناوری را ارائه خواهیم داد که برای ارزیابی در هر یک از موارد می تواند به کار برده شود. چون خزش گر، برای سرور وب یک سرویس گیرنده مانند دیگران است از اصطلاح «بایگانی جانبی سرویس گیرنده» برای این روش فراهم آوری استفاده می کنیم. بسته به ساختار پسین وب و سطح تعامل با سرویس گیرنده خزش گر ها می توانند یا یک وبگاه کامل و یا بخشی از اجزای آن را به تنهایی تصرف کنند. جزء باقی مانده برای خزش گرها غیر قابل دسترسی است در اصطلاح شناسی موتور جست و جو «وب عمیق» یا وب پنهان نامیده شده است.

این اصطلاح شناسی را تصدیق خواهیم کرد تا زمانی که مشخص شود که تعیین حدود وب مخفی، به طور محض فنی بوده و به طور دائم به عنوان خزنده ها توانایی شان را در یافتن راهی برای اسناد، بهبود بخشید. دو روش متناوب برای گردآوری محتوا باقی می ماند، حتی اگر بسیار کم به کار برده شوند و حتی هم چنان تحت تحقیق و بررسی مانده باشند هر دو نیاز به فعالیت از جانب سرور دارند که نه تنها به یک مجوز نیاز دارد بلکه همچنین به یک مشارکت فعال انتشار دهنده سایت برای استفاده شدن نیز نیاز دارد. اولی، بر اساس کاربران سایت است که مسیر هدایت را توسعه می دهد و شناسایی محتوای سایت را برای بایگانی آن انجام می دهد و چون بر اساس ثبت تراکنش های ایجاد شده بین کاربران سایت و سرور می باشد آن را بایگانی تراکنشی می نامیم.

دومی، شامل بایگانی به طور مستقیم از ناشر اجزای گوناگون سیستم اطلاعات وب و انتقال آن ها به یک شکل بایگانی شده می باشد بنابراین بایگانی جانبی سرور نامیده می شود. این روش های متناوب، سخت تر از بایگانی جانبی سرویس گیرنده می باشند چون همانطور که در بالا ذکر شد، نه تنها به یک مشارکت فعال از جانب سازندگان نیاز دارند بلکه باید بر مبنای مورد به مورد اجرا گردند اما حتی اگر افزایش مقیاس نداشته باشند می توانند در مواردی به کار برده شوند؛ برای مثال در جایی که خزنده برای تصرف دقیق موفق نمی شود و زمانی که محتوا در آن کمتر به کار رفته باشند.

4-2-1- بایگانی جانبی سرویس گیرنده

این روش مهمی در فراهم آوری است در هر دو، به علت سادگی، قابلیت مقیاس پذیری، و تطبیق با یک محیط سرور - سرویس گیرنده می باشد (شکل 1-1 را ببینید). خزنده ها با آن چه که روش معمول در دستیابی به وب است سازگار شده اند و این کار بایگانی هر سایت را مجاز می سازد که آزادانه به وب باز شبکه داخلی یا خارجی دسترسی داشته باشند تا زمانی که خزنده به یک اجازه مناسب دست یابد.



تصویر ۱- آرشیو سمت کاربر (Client-side): آرشیو وب در تقابل با آرشیو کاربر است تا محتوا را از سرور وب جمع‌آوری کند. سرور وب می‌تواند محتوا را از سرورهای متعدد و متنوعی فراهم کند (درخواست‌ها، پایگاه داده‌ها، فایل‌های سرور).

این روش، نه تنها موقعیت مشابهی را که کاربران عادی وب نیاز دارند، اتخاذ می‌کند؛ بلکه شکل تعامل‌هایش با سرورها را تقلید می‌کند. خزش‌گرها از صفحه‌ها هسته شروع می‌شوند، آنها را تجزیه می‌کنند، پیوندها را برداشت می‌کنند، و سند پیوندی را واکنشی می‌کنند؛ سپس آنها این پردازش با سند واکنشی شده را تکرار می‌نمایند و تا زمانی که پیوندهایی را کشف کنند^۱ و سند را درون حوزه تعیین شده پیدا کنند. این پردازش مورد نیاز است، چون HTTP فرماتی را که باید فهرست کامل سند قابل دسترس بر روی سرور را باید بازگرداند که برای مثال بر خلاف FTP است. از این رو، هر صفحه باید به وسیله پیوندی از صفحه‌ها دیگر «کشف» شود.

فناوری خزش، برای اهداف نمایه‌سازی توسعه یافته است^۲. به کار بردن آن برای بایگانی وب، برخلاف اینکه جنبه‌های بیشتری از این فناوری را دوباره استفاده می‌کند، تغییراتی را برای آن ایجاد می‌کند.

نخست اینکه، خزش‌گرهای بایگانی باید برای واکنشی تمام فایل‌ها تلاش کنند هرچه فرمتشان برای بایگانی یک مدل کامل از سایت‌ها باشد بر خلاف خزش‌گرهای موتور کاوش که معمولاً فقط فایل‌هایی را واکنشی می‌کند که آنها بتوانند فهرست نمایند. خزش‌گرهای موتور کاوش، برای مثال اغلب از انتقال از

۱. برای نظرات اخیر درباره فناوری خزشگر پلنت و دیگران (۲۰۰۴) و چاکرلبارتی (۲۰۰۲) را ببینید.
۲. برای نظر اخیر درباره توسعه موتور کاوش تجاری سونتریج (۱۹۹۷) را ببینید.

تصویر ۱- آرشیو سمت کاربر (client-side): آرشیو وب در تقابل با آرشیو کاربر است تا محتوا را از سرور وب جمع‌آوری کند. سرور وب می‌تواند محتوا را از سرورهای متعدد و متنوعی فراهم کند (درخواست‌ها، پایگاه داده‌ها، فایل‌های سرور).

این روش، نه تنها موقعیت مشابهی را که کاربران عادی وب نیاز دارند، اتخاذ می‌کند؛ بلکه شکل تعامل‌هایش با سرورها را تقلید می‌کند. خزش‌گرها از صفحه‌ها هسته شروع می‌شوند، آنها را تجزیه می‌کنند، پیوندها را برداشت می‌کنند، و سند پیوندی را واکنشی می‌کنند؛

سپس آن‌ها این پردازش با سند واکنشی شده را تکرار می‌نمایند و تا زمانی که پیوندهایی را کشف کنند (1) و سند را درون حوزه تعیین شده پیدا کنند این پردازش مورد نیاز است چون HTTP فرمانی را که باید فهرست کامل سند قابل دسترس بر روی سرور را باید بازگرداند که برای مثال بر خلاف FTP است. از این رو هر صفحه باید به وسیله پیوندی از صفحه‌ها دیگر «کشف» شود.

فناوری خزش، برای اهداف نمایه‌سازی توسعه یافته است (2). به کار بردن آن برای بایگانی وب، بر خلاف اینکه جنبه‌های بیشتری از این فناوری را دوباره استفاده می‌کند تغییراتی را برای آن ایجاد می‌کند.

نخست این که خزش گره‌های بایگانی باید برای واکنشی تمام فایل‌ها تلاش کنند هر چه فرمت‌شان برای بایگانی یک مدل کامل از سایت‌ها باشد بر خلاف خزش گره‌های موتور کاوش که معمولاً فقط فایل‌هایی را واکنشی می‌کند که آن‌ها بتوانند فهرست نمایند. خزش گره‌های موتور کاوش برای مثال اغلب از انتقال از

ص: 28

1- برای نظرات اخیر درباره فناوری خزشگر پلنت و دیگران (2004) و چاکرابارتی (2002) را ببینید.

2- برای نظر اخیر درباره توسعه موتور کاوش تجاری سونریج (1997) را ببینید.

فایل های کاربردی و ویدئویی بزرگ چشم پوشی می کند. بارگذاری این نوع فایل ها می تواند تفاوت قابل توجهی را در رابطه با زمان و پهنای باند مورد نیاز برای خزش در کل سایت ایجاد نماید.

تفاوت دوم، با مدیریت موقتی خزش ها مرتبط می باشد. برای اجتناب از اضافه بار سرورهای وب، خزش گر ها از قوانین مطلوب و معتبری استفاده می کنند (روشه، 2006). این مستلزم این است که تصرف یک وب می تواند در طی چندین دقیقه منتهای مراتب چندین ساعت و گاهی اوقات چندین روز طول بکشد. یک محاسبه ساده نشان می دهد زمانی را که در خصوص یک تأخیر 3 ثانیه ای بین دو درخواست، بیشتر از 3 روز برای بایگانی یک سایت با 100/000 صفحه طول خواهد کشید. این تأخیر، مبحث ثبات موقت تصرف را که دستخوش تغییراتی در طول زمان می شود، افزایش می دهد. اگر صفحه نمایه برای مثال در طی تصرف تغییر یابد روش بایگانی با آخرین بایگانی که با صفحه ها بایگانی پیوند داده شده است سازگار نمی باشد.

این یک مورد برای خزش گر های بایگانی است چون خزشگر برای تهیه محتوا فرض شده است البته نه فقط در جهت هدایت محتواها خزش گر های موتور کاوش عادت دارند که به صفحه ها زنده اشاره کنند. به این معنی که مضمون ابرمتن برای آن ها، یکی است که توسط سرور اصلی تهیه شده است (که البته در سراسر صفحه ها دارای ثبات است و به هنگام سازی شده است). بر عکس، خزش گر های بایگانی باید محتوا را به طور کلی تصرف کنند که انجام خواهند داد با وابستگی یا بدون وابستگی داخلی، به عنوان فقط مضمون هایی برای شناوری و تفسیر.

این امر به مراتب نتایج قابل دسترسی نسبت به خط مشی خزش گرها دارد. چون مطلوبیت برای سرورها حکم می کند که همیشه عملیات خیلی محدودی برای خزش داشته است، خزش گر های موتور های کاوش از اولویت خزش در سطح پهنای نخست استفاده کرده اند با برخی تغییرات اساسی که با خزش در بهترین صفحه های اول، هدفمند شده است (چو و همکاران، 1998؛ ناجون و هیدان 2001) (1)؛ ناجورک و واینر 2001 (2)؛ کاستیلو و همکاران، 2004؛ بازا - یس و کاستیلو (2005) (3).

اتخاذ این خط مشی، همچنین، روش برای به حداقل رسانیدن اثرات شدید تله های روباتیک (4) بر روی کل خزش با قرار گیری بیرون از خزش بر روی تعداد زیادی از سایت های مختلف می باشد.

اما این راهبرد زمان بندی خزش، در دسر افزایش تفاوت موقتی خزش ها در سطح سایت را دارد. بنابراین برای اتخاذ خزش های بایگانی در اولویت نخست یک سایت پیشنهاد شده است. (5) اما برای خزش های مقیاس بزرگ، هنوز برای بهینه سازی بازده خزننده با اطمینان از منابعی که در حداکثر ظرفیت خود استفاده شده اند، ضروری است تأخیر آشکار بین درخواست ها و منابع قابل دسترس خزش

ص: 29

Najork and Heydon -1

Wiener -2

Baeza-Yates and Castillo -3

4- فاصله ای عمودی در یک سیستم پردازش اطلاعات که به منظور جمع آوری تغییر یا خراب کردن آتی اطلاعات به وجود آمده است
5- به طور موردی، برای ملزومات درون بحث شده است (میسانس، 2004)؛ در مورد سیاست های الگوریتم زمان بندی خزش که در سایت به عنوان بعدی عمودی آمیخته می شود منابع زیر را ببینید (Castillo et al. 2004, Baeza-Yates and Castillo (2005)).

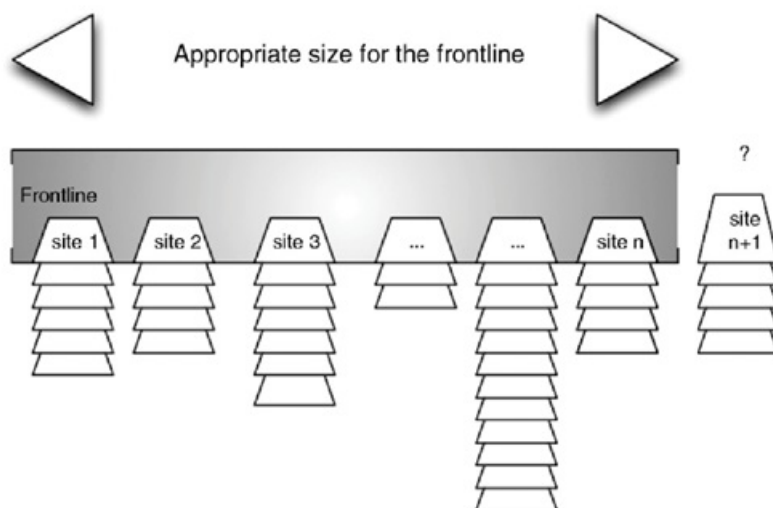
موجود است، شخص باید تعداد مطلوبی از سایت‌ها را برای شروع در زمان مشابه با ایجاد فرکانس‌های مطمئن جست و جو کند که به وسیله قوانین مطلوب بدون هیچ تأخیر غیر ضروری بین درخواست‌ها واقع شده‌اند. شکل 2-1 خط مقدم یک خزش‌گر را نشان می‌دهد که اندازه‌ای مناسب با

تخصیص بهینه منابع خزش دارد.

عکس

۳۰ مدیریت منابع اطلاعاتی وب

موجود است، شخص باید تعداد مطلوبی از سایت‌ها را برای شروع در زمان مشابه با ایجاد فرکانس‌های درخواست‌های مطمئن جست‌وجو کند که به وسیله قوانین مطلوب، بدون هیچ تأخیر غیر ضروری بین درخواست‌ها واقع شده‌اند. شکل 2-1 خط مقدم یک خزش‌گر را نشان می‌دهد که اندازه‌ای مناسب با تخصیص بهینه منابع خزش دارد.



تصویر 2-1. Frontline شامل سایت‌هایی است که باید با یک خزشگر یکسان و به‌طور همزمان خزش شوند. اندازه آن (n) در حد اَپتیمم است و اگر بین درخواست‌ها تأخیری رخ دهد، فقط به‌وسیله قوانین ساده‌ای محدود می‌شود و منابع خزش همچنان درگیر و مشغول خواهند بود. اگر $n+1$ سایت خزش شوند، محدودیت منابع خزش نوع تأخیر اضافی و عدم ربط موقت را به کار خواهند گرفت. اگر $n-1$ سایت خزش شوند، منابع بدون استفاده خواهند ماند.

محدودیت‌هایی برای هر آنچه که بتوان با استفاده از این روش بایگانی کرد وجود دارد. بیش از همه طی برداشت پیوند و برخی در طی بازیابی از طریق واسط HTTP رخ می‌دهد. دلیل مورد قبلی این حقیقت است که URL استخراج شده، به‌طور نامساعد شکل گرفته است یا از پارامترهای پیچیده استفاده کرده‌اند، یا به‌سختی برای تجزیه کردن URL از فایل آغازگر یا فایل اجرا یا حتی کد HTML استفاده کرده است. دومی می‌تواند به‌علت تجدید مسیرها، مذاکره محتوا، اجازه، پاسخ‌های تدریجی (کند)، اندازه نهایی، اتصالات TCP استثناهایی، پاسخ‌های سرور نامعتبر، و مانند آن باشد. با استفاده از این نوع ابزارها، فراهم‌آوری مقیاس بزرگی از متوا در یک مسیر کل نگ، که البته این طور نیست، مجاز می‌شود.

تصویر 1-2 Frontline شامل سایت هایی است که باید با یک خزشگر یکسان و به طور همزمان خزش شوند. اندازه آن (n) در حد اپتیمم است و اگر بین درخواست ها تأخیری رخ دهد فقط به وسیله قوانین ساده ای محدود می شود و منابع خزش همچنان درگیر و مشغول خواهند بود. اگر n+1 سایت خزش شوند، محدودیت منابع خزش نوع تأخیر اضافی و عدم ربط موقت را به کار خواهند گرفت. اگر n-1 سایت خزش شوند، منابع بدون استفاده خواهند ماند.

محدودیت هایی برای هر آن چه که بتوان با استفاده از این روش بایگانی کرد وجود دارد. بیش از همه طی برداشت پیوند و برخی در طی بازیابی از طریق واسط HTTP رخ می دهد. دلیل مورد قبلی این حقیقت است که URL استخراج شده به طور نامساعد شکل گرفته است یا از پارامترهای پیچیده استفاده کرده اند یا به سختی برای تجزیه کردن URL از فایل آغازگر یا فایل اجرا یا حتی کد HTML استفاده کرده است. دومی می تواند به علت تجدید مسیرهها، مذاکره محتوا، اجازه، پاسخ های تدریجی (کند)، اندازه نهایی، اتصالات TCP استثنایایی، پاسخ های سرور نامعتبر، و مانند آن باشد.

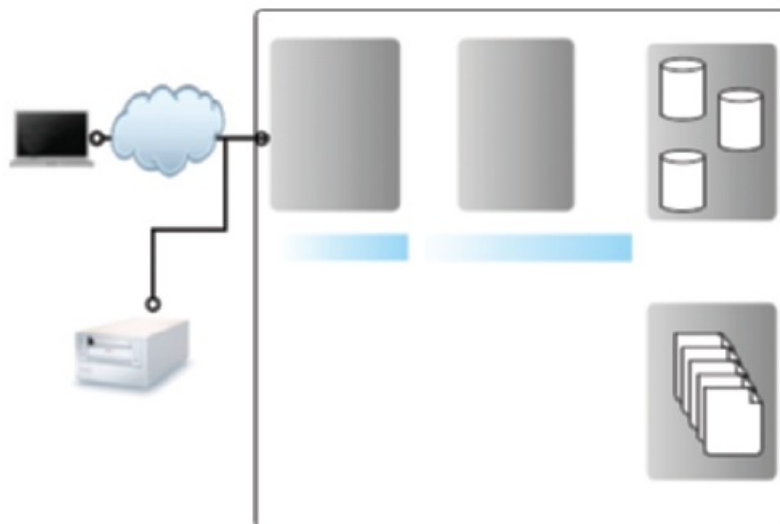
با استفاده از این نوع ابزارها، فراهمآوری مقیاس بزرگی از متوا در یک مسیر کل نگ، که البته این طور نیست. مجاز می شود.

بایگانی تراکنشی (شکل 3-1 را ببینید) به وسیله فیچ (2003) (1) پیشنهاد شده است که شامل تصرف و بایگانی «تمام پاسخ های متمایز اساسی که توسط یک وبگاه تولید شده است و در رابطه با محتوایشان و چگونگی تولید می باشد». این کار در سیستم پرش صفحه (2) با استفاده از یک فیلتر درون ورودی سرور وب (درخواست) جریان و جریان خروجی (پاسخ) اجرا شده است. این توابع عملیاتی اکنون بر روی برخی از نظام های مدیریت محتوای وب مانند Vignette TM قابل دسترس هستند.

عکس

۴-۲-۲- بایگانی تراکنشی

بایگانی تراکنشی (شکل ۳-۱ را ببینید) به وسیله فیچ (۲۰۰۳) پیشنهاد شده است که شامل تصرف و بایگانی «تمام پاسخ‌های متمایز اساسی که توسط یک وبگاه تولید شده است و در رابطه با محتوایشان و چگونگی تولید می‌باشد». این کار در سیستم پرش صفحه^۱ با استفاده از یک فیلتر درون ورودی سرور وب (درخواست) جریان و جریان خروجی (پاسخ) اجرا شده است. این توابع عملیاتی اکنون بر روی برخی از نظام‌های مدیریت محتوای وب مانند Vignette TM قابل دسترس هستند.



شکل ۳-۱ بایگانی تراکنشی

جفت‌های درخواست/ پاسخ منحصر به فرد ذخیره و بایگانی شده‌اند از این رو، ایجاد یک بایگانی کامل از تمام محتوا برای یک سایت مشخص پیش‌بینی شده است. درخواست‌ها فقط با اندکی تفاوت (غیر مادی) به صورت منحصر به فرد مورد بررسی قرار می‌گیرند، به استثنای محاسبه مجموع مقابله‌ای قسمتی از کد، که آنها را به صورت رمز در آورده است، کیفیت دقیق اینها چگونه می‌تواند با تعداد زیادی از روش‌های شخصی‌سازی محتوا، که واضح نمی‌باشد، منطبق گردد. این نوع بایگانی وب می‌تواند فواید پیگردی و ثبت هر برنامه‌ریزی ممکن در محتوا را ثابت کند.

1. Fitch
2. <http://www.projectComputing.com/products/pageVault>

شکل 3-1 بایگانی تراکنشی

جفت‌های درخواست/ پاسخ منحصر به فرد ذخیره و بایگانی شده‌اند از این رو، ایجاد یک بایگانی کامل از تمام محتوا برای یک سایت مشخص پیش‌بینی شده است. درخواست‌ها فقط با اندکی تفاوت («غیر مادی») به صورت منحصر به فرد مورد بررسی قرار می‌گیرند، به استثنای محاسبه مجموع مقابله‌ای قسمتی از کد، که آن‌ها را به صورت رمز در آورده است،

کیفیت دقیق اینها چگونه می‌تواند با تعداد زیادی از روش‌های شخصی‌سازی محتوا، که واضح نمی‌باشد، منطبق گردد.

این نوع بایگانی وب می تواند فواید پیگردی و ثبت هر برنامه ریزی ممکن در محتوا را ثابت کند.

ص: 31

Fitch -1

<http://www.projectComputing.com/products/page Vault> -2

محتوایی که هرگز دیده نشده، بایگانی نخواهد شد (همان طور که ذکر شد، بوفخواد و وینناد (2003) (1)، برآورد کردند که 25 درصد صفحه ها از یک وبگاه بزرگ علمی، هرگز قابل دستیابی نیستند). اما محتوای وب پنهان تا زمانی که به دستیابی برسد ثبت خواهد شد و این یک مزیت مهم است.

محدودیت مهم این روش این است که باید با موافقت و همکاری مالک، سرور، اجرا شود. از این رو، برای بایگانی داخلی وب نشان داده شده است. و این مسئله مزیت توانا بودن ثبت دقیق از چیزی را که و زمانی را که دیده است دارا می باشد. برای بایگانی شرکت یا سازمانی، اغلب به وسیله جواب گویی قانونی، برانگیخته می شود. این کار، می تواند یک مزیت باشد. حتی برای ترکیب با اطلاعات از سرور ثبت وقایع درباره کسی که محتوا را دیده است، امکان پذیر باشد. به طور واضح، هر چه را که می تواند به عنوان یک مزیت برای بایگانی وب داخلی دیده شود ممکن است برای یک بایگانی عمومی مشکل باشد، چون می تواند تأکید بر محرمانگی جدی را بالا ببرد. اما به هر حال محتوای قابل استفاده نمی باشد.

3-2-4-بایگانی سرور - جانبی

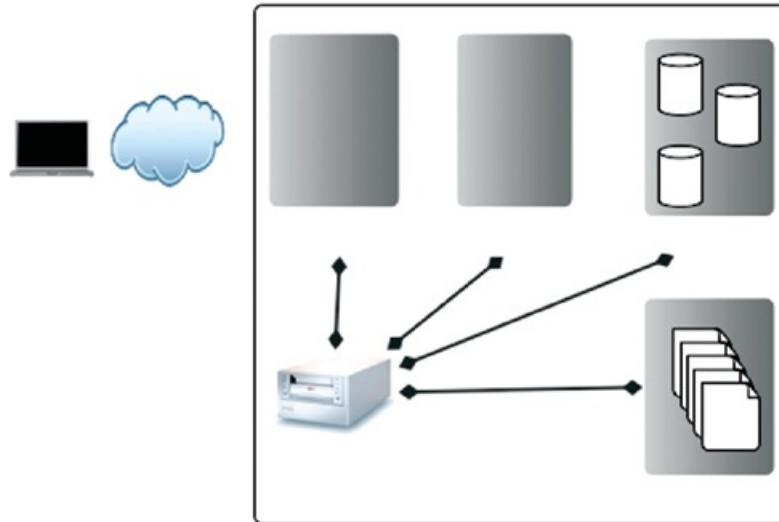
نوع آخر از روش فراهم آوری بایگانی وب، کپی مستقیم فایل ها از سرور، بدون استفاده از واسطه های HTTP است. این روش مانند روشی قبلی می تواند فقط با همکاری مالکان سایت استفاده گردد (شکل 4-1).

عکس

محتوایی که هرگز دیده نشده، بایگانی نخواهد شد (همان‌طور که ذکر شد، بوفخواد و وینناد^۱، ۲۰۰۳ برآورد کردند که ۲۵ درصد صفحه‌ها از یک وبگاه بزرگ علمی، هرگز قابل دستیابی نیستند). اما محتوای وب پنهان تا زمانی که به دستیابی برسد، ثبت خواهد شد و این یک مزیت مهم است. محدودیت مهم این روش، این است که باید با موافقت و همکاری مالک سرور، اجرا شود. از این رو، برای بایگانی داخلی وب نشان داده شده است. و این مسئله مزیت توانا بودن ثبت دقیق از چیزی را که و زمانی را که دیده است دارا می‌باشد. برای بایگانی شرکت یا سازمانی، اغلب به وسیله جوابگویی قانونی، برانگیخته می‌شود. این کار، می‌تواند یک مزیت باشد. حتی برای ترکیب با اطلاعات از سرور ثبت وقایع درباره کسی که محتوا را دیده است، امکان‌پذیر باشد. به‌طور واضح، هر چه را که می‌تواند به‌عنوان یک مزیت برای بایگانی وب داخلی دیده شود، ممکن است برای یک بایگانی عمومی مشکل باشد، چون می‌تواند تأکید بر محرمانگی جدی را بالا ببرد. اما به هر حال محتوا قابل استفاده نمی‌باشد.

۴-۲-۳-بایگانی سرور- جانبی

نوع آخر از روش فراهم آوری بایگانی وب، کپی مستقیم فایل‌ها از سرور، بدون استفاده از واسطه‌های HTTP است. این روش مانند روشی قبلی می‌تواند فقط با همکاری مالکان سایت استفاده گردد (شکل ۴-۱).



شکل ۴-۱. آرشیو Server - side: قطعه‌های مختلف اطلاعات به‌طور مستقیم از سرورها آرشیو می‌شوند. تولید نمونه کار محتوای آرشیو شده و یک نمونه پشتیبان (backup) از فایل‌ها چالش اصلی این روش است.

1. Boufkhad and Viennot

اگر چه بسیار ساده به نظر می رسد، در واقع مشکل زیادی را برای ایجاد محتوای کپی شده قابل استفاده افزایش می دهد و حتی در مورد فایل های ایستای HTML، شخص می تواند به زحمت در محتوا از طریق پیوندهای مطلق به عنوان نام دامنه هدایت شود که در بایگانی متفاوت خواهند بود. اما بیشترین مشکل ناشی از محتوای توسعه یافته پویاست که تکه های به هم پیوسته محتوا از منابع گوناگون (قالب ها (1) و پایگاه های داده) است که توسط درخواست های کاربرد در حالت پرواز در ارتفاع کم توسعه یافته است. کپی برداری فایل های پایگاه داده ها، قالب ها، و فایل های آغاز گر به این معنی نیست که آن برای تولید مجدد محتوا از بایگانی آسان خواهد بود. بر عکس یک وظیفه چالش برانگیز خواهد بود چون نیازمند توانا شدن برای اجرای در محیط مشابه، با پارامترهای مشابه در بایگانی است. در واقع، وقتی امکان پذیر شد، محتوای توسعه یافته پویا در شکل نهایی اش بهتر حفظ می شود، معمولاً در فایل های HTML مسطح بهتر حفظ می شود (برای مثال این موردی برای بیشتر CMSها، بلاگ ها و ویکی ها (2) صورت می گیرد).

اما گاهی اوقات این کار مشکل است، حتی برای خزش گرهای یافتن مسیر برای برخی اسناد یک وبگاه و فایل هایی که می تواند از طریق یک تعامل پیچیده به دستیابی برسند، غیر ممکن است (مانند ارسال یک پرس و جو و یک فرم یا پنجره) که به سختی توسط خزش گرها تصرف خواهد شد این بخش وب، وب پنهان یا عمیق نامیده شده است (برگمان، 2001؛ چانگ و همکاران 2004) که بزرگ تر از وب سطحی است و (همچنین به طور عمومی وب قابل نامیده شده است). (3)

در این مورد، بایگانی جانبی سرور می تواند یک راه حل باشد، همان طور که در بالا ذکر شد، به مشارکت فعال مجری سایت نیاز دارد. بیشتر از یک پشتیبانی ساده است که دسترسی به محتوا را در نمونه های اصلی اش تضمین نمی کند، آن بر توانا بودن برای «نمایش» مجدد سایت در محیط بایگانی را دلالت می کند. این کاهش وابستگی به پایگاه داده ها و اجرای فایل های آغاز گر جانبی سرور به قدری که امکان پذیر است را نشان می دهد. این کار می تواند به وسیله استخراج اطلاعات ساخت یافته محتوی در پایگاه داده ها و انتقال آن درون XML انجام شود. نوعی معماری اطلاعات دروازه اسناد نامیده شده است که شامل اسناد غیر وب با آن هایی که به وسیله کاتالوگی می تواند مانند اینها بایگانی شود، قابل دستیابی می باشد. این کار، برای چندین سایت انجام شده است که به مقوله سایت های پنهان به وسیله کتابشناسی ملی فرانسه (4) وابسته می باشد (فصل 5 را ببینید).

این کار فقط در چارچوب واسپاری قانونی امکان پذیر می باشد که در فرانسه مانند بسیاری از کشورهای دیگر به کار گرفته می شود. حقیقت این است که وب پنهان همچنین غالب اوقات دارای محتوای بسیار غنی با این نوع معماری اطلاعات است که انبوه زیادی از محتوای از قبل موجود بر روی وب منتشر کرده است. عمومیت این نوع معماری اطلاعات، بایگانی سایت سرور را می سازد، روشی که به جایی که می تواند به کار رود توجه دارد.

ص: 33

templatestemplates -1

wikis -2

3- این اصطلاح برای قسمتی از وب مشخص شده است که می تواند توسط خزش گرهای نمایه شود (لورنس و گیلز، 1998، 1999)

Bibliothèque nationale de France -4

همان طور که قبلاً دیده شد کپی برداری از یک وبگاه، یک وظیفه غیر پیش پا افتاده است. آن در واقع بر ایجاد مجدد یک سیستم اطلاعاتی اشاره می کند که برای کاربران قابل دسترسی خواهد بود. همان طور که آنتونیول و همکارانش (1) آن را در وبگاه قرار دادند «وبگاه ممکن است به سادگی یک فایل واحد یا یکی از پیچیده ترین مجموعه های محصولات مصنوعی نرم افزاری مشترک باشد که تا به حال درک شده است». به طور مطلوب، بایگانی می توانست در اصل متناظر (هم ریخت) ساختار سلسله مراتبی مشابه، نام گذاری فایل ها، سازوکار پیونددهی، چارچوب) باشد، اما به دلایل عملی اینگونه نیست. همان طور که در بخش قبلی دیده شده اکتساب سایت ها در حقیقت موارد یک تغییر شکل فایل را به طور مؤثر کاهش می دهد.

چالش بیشتر، در ایجاد مجدد سیستم های اطلاعاتی مشابه می باشد. نظام اطلاعاتی وب معماری اطلاعاتی پیچیده را نشان می دهد که به سیستم های عملیاتی خاص پیکربندهای سرور و محیط کاربردی وابسته هستند که در بیشتر موارد حتی برای ایجاد مجدد استفاده از (فایل) چرکنویس برای طراحان و مدیران مشکل است به همین دلیل است آرشیویست ها مجبور به اتخاذ راهبرد تبدیل یا تغییر شکل می شوند. این تبدیل ها می توانند اثر شدیدی بر ساز و کارهای آدرس دهی و پیوندی، فرمت ها، و همچنین تغییر خود شیء تأثیر داشته باشند.

تاکنون سه راهبرد برای ساخت بایگانی وب اتخاذ شده است. راهبرد نخست برای ایجاد کپی محلی از فایل های سایت و هدایت از طریق این کپی در یک مشابه برای مثال بر روی وب می باشد.

راهبرد دوم، برای اجرای یک سرور وب با به کار بردن مضمون در یک محیط برای مرورگرهای کاربران است راهبرد دوم راه اندازی سرور وب و به کارگیری محتوا در این محیط برای مرورگرهای کاربران است. سومین راهبرد، سازماندهی مجدد اسناد منطبق بر نام گذاری متفاوت (غیر وب)، آدرس دهی منتقل کردن است. بخش های زیر موافقان و مخالفان این رهنمودهای متفاوت و همچنین موارد استفاده را نشان می دهد.

1-3-4- نظام فایل های محلی به خدمت گرفته شده بایگانی ها

توصیف

این نوع بایگانی (شکل 5-1) بر اساس احتمالاتی است که مشخصه های URL از پیشوند سیستم فایل محلی، «فایل» (2) استفاده کند در یک شمای URL برای کپی و دسترسی به فایل های محلی از وبگاه اصلی استفاده کند مانند:

[Http://www.example.org/example.HTML](http://www.example.org/example.HTML)

[File:///Users/archire2005/eample.org/example.HTML](file:///Users/archire2005/eample.org/example.HTML)

این، استفاده از سیستم فایل محلی را برای شناسایی از طریق مواد بایگانی شده در وب، قادر می سازد.

همچنین به استفاده از شکل جزئی (وابسته) از URL که نه تنها از پیشوند دوری می کند بلکه همچنین از نام سرور و مسیر هدف نیز اجتناب می کند.

مرورگرهای استاندارد می توانند به طور مستقیم باز شوند (به طور مثال بدون سرور وب) و چنین فایل های ذخیره شده را محلی نمایند و تا زمانی که پیوندها اسناد وابسته هستند، شناوری در بایگانی مشابه مانند همانی است که در سایت اصلی خواهد بود، به جز فقط در قسمت آدرس دهی (1) مرورگر وقتی در جست و جوی پیشوند URL هستند. (در اینجا «فایل» به جای HTTP می باشد).

توضیح

فایده اصلی این راهبرد، تسهیل دستیابی به بایگانی با مسیره‌ی ساختار و بگانه اصلی به سوی نظام بایگانی فایل است. با استفاده از مرورگر استاندارد و سیستم فایل مجاز به اجتناب از هزینه های بالا که همراه است با اجرای دسترسی به بایگانی مبتنی بر سرور. بنابراین، حتی گروهی با مهارت های فنی فناورانه بسیار اساسی می توانند این نوع بایگانی را بر پا و اجرا کنند. ولی محدودیت هایی در این رویکرد وجود دارد. از یک نگرش محافظه کارانه نقص اصلی این است که تبدیل های گوناگون فایل های اصلی، مورد نیاز می باشند. از این رو، صداقت محض برای فایل های اصلی، نمی تواند مورد ملاحظه قرار گیرد مگر به وسیله سندسازی با تغییرات دقیق که در فایل اصلی به کار می رود یا به وسیله نگهداری یک کپی از اصل. انتقال محتوا در دو سطح در رویکرد بایگانی «FS» محلی مورد نیاز می باشد.

نخست، به علت تفاوت در نام قرار داد بین URL و نظام فایل (تصویب و ذخیره کاراکترها، رهایی از قوانین، حساسیت موردی)، نام اشیاء ممکن است به تغییراتی نیاز داشته باشد (فصل 1 - B را برای نشان دادن جزئیات بیشتر این تغییرات ببینید). در مورد جایی که صفحه طبق پارامترها پرس و جو شده و به طور چشمگیری تولید شده است، یک نام حتی برای ایجاد شدن در صفحه نتیجه، پارامترهایی را برای اطمینان از منحصر به فرد بودن صفحه بایگانی شده تأمین می کند.

دوم اینکه، پیوندهایی مطلق باید در پیوندهای مرتبط در خودکد صفحه منتقل شوند تا شناوری مبتنی بر نظام فایل را امکان پذیر کنند. حتی اگر این کار بتواند به سادگی توسط تغییر شکل URL اصلی درون یک توضیح در کد، سندیت داده شود این دستکاری منبع را نشان می دهد.

از نظر عملی، اشکالات اصلی ناشی از خود نظام فایل می باشد که یک معماری اطلاعات متفاوت قابل توجه نسبت به وب دارد. نخست اینکه، سازماندهی بایگانی باید در سازماندهی سلسله مراتبی نظام های فایل مناسب باشد، با این حال، یک بایگانی نه تنها از ترکیب سایت ها، بلکه همچنین گروه هایی در سایت ها (مجموعه ها) و نسخه هایی از سایت ها تشکیل شده است. مسیره‌ی این سازماندهی برای یک ساختار سلسله مراتبی، بدون تغییر و انتخاب روی نمی دهد. چگونه سایت ها در یک حالتی که زمان پایدار، یک مبحث مهم برای بررسی می باشد، با یکدیگر در یک گروه قرار می گیرند. نام های مجموعه

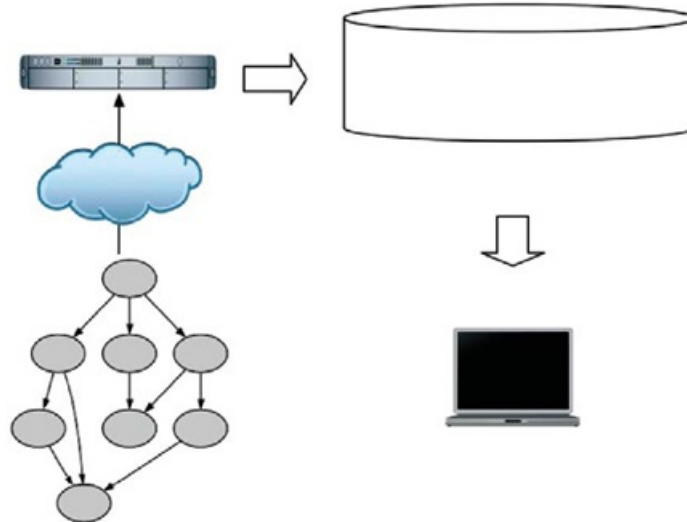
باید دائمی باشند و زمان گروه بندی نیز باید برای ضبط فرکانس ها، تطبیق یافته باشد. در مسیر تمام این مباحث، تصمیم گیری کامل باید پیشاپیش انجام شده باشد. آن ها بر روی چگونگی انتخاب ساختار که در برابر توسعه مجموعه مقاومت می کنند، تحت فشار قرار می گیرند. سازماندهی زمان در اجرای هر جمله از یک برنامه برای اهداف اشکال زدایی (از یک مدل به مدل دیگر) سایت نیز یک مبحث مهم برای هر یک از لایه های نرم افزار است که باید به سیستم فایل با استاندارد بالا افزوده شود. این لایه باید قادر باشد حداقل مدل های متفاوت سایت ها را که به تاریخ شان بستگی دارد با یکدیگر متصل نماید (نسخه پردازی) و این کار را با یک واسط کاربر مناسب برای شناسایی از طریق زمان به سادگی ارائه می دهد. این کار اغلب با استفاده از مدیریت خارجی یک پایگاه داده در سایت ها و ضبط اطلاعات و ابزارهایی برای ایجاد صفحه های نمایش میانجی با فهرستی از تاریخ اجرا می شود که در آن سند بایگانی شده است.

محدودیت دیگر این رویکرد، ناشی تعداد عظیم فایل های بایگانی وب است که باید جا به جا شوند. این مسئله عادی است که بایگانی هایی با بیلیون ها صفحه را ببینیم. این شکل به محدودیت های ظرفیت نظام های فایل جاری می رسد و دسترسی می یابد حتی وقتی یک نظام فایل بتواند این مقدار فایل را جا به جا کند، عملکرد می تواند تحت تأثیر قرار گیرد. برای کم کردن باری که بر روی نظام فایل قرار گرفته، بایگانی با مقیاس بزرگ برای فایل های مخزن استفاده شده است. البته این کار، ارتباط مستقیم در نام گذاری و پیوند دادن را که رویکرد بایگانی نظام فایل محلی برای اتخاذ رویکرد دوم ارائه می دهد را می شکند و بایگانی مبتنی بر سرور به کار رفته در وب برای تحویل محتوا از این فایل مخزن است.

شکل 1-5. بایگانی سیستم فایل محلی. فایل اصلی مورد خزش قرار می گیرد و صفحه ها و سایر فایل ها به صورت انفرادی روی سیستم فایل بایگانی ذخیره می شوند. دست یابی بوسیله شناسایی مستقیم در نظام فایل انجام می شود.

عکس

باید دائمی باشند و زمان گروه‌بندی نیز باید برای ضبط فرکانس‌ها، تطبیق یافته باشد. در مسیر تمام این مباحث، تصمیم‌گیری کامل باید پیشاپیش انجام شده باشد. آنها بر روی چگونگی انتخاب ساختار که در برابر توسعه مجموعه مقاومت می‌کنند، تحت فشار قرار می‌گیرند. سازماندهی زمان در اجرای هر جمله از یک برنامه برای اهداف اشکال‌زدایی (از یک مدل به مدل دیگر سایت) نیز یک مبحث مهم برای هر یک از لایه‌های نرم‌افزار است که باید به سیستم فایل با استاندارد بالا افزوده شود. این لایه باید قادر باشد حداقل مدل‌های متفاوت سایت‌ها را که به تاریخ‌شان بستگی دارد با یکدیگر متصل نماید (نسخه‌برداری) و این کار را با یک واسط کاربر مناسب برای شناسایی از طریق زمان به سادگی ارائه می‌دهد. این کار اغلب با استفاده از مدیریت خارجی یک پایگاه داده در سایت‌ها و ضبط اطلاعات و ابزارهایی برای ایجاد صفحه‌های نمایش میانجی با فهرستی از تاریخ اجرا می‌شود که در آن سند بایگانی شده است. محدودیت دیگر این رویکرد، ناشی تعداد عظیم فایل‌های بایگانی وب است که باید جابه‌جا شوند. این مسئله عادی است که بایگانی‌هایی با بیلیون‌ها صفحه را ببینیم. این شکل به محدودیت‌های ظرفیت نظام‌های فایل جاری می‌رسد و دسترسی می‌باید حتی وقتی یک نظام فایل بتواند این مقدار فایل را جابه‌جا کند، عملکرد می‌تواند تحت تأثیر قرار گیرد. برای کم کردن باری که بر روی نظام فایل قرار گرفته، بایگانی با مقیاس بزرگ برای فایل‌های مخزن استفاده شده است. البته این کار، ارتباط مستقیم در نام‌گذاری و پیوند دادن را که رویکرد بایگانی نظام فایل محلی برای اتخاذ رویکرد دوم ارائه می‌دهد را می‌شکند و بایگانی مبتنی بر سرور به کار رفته در وب برای تحویل محتوا از این فایل مخزن است.



شکل ۱-۵. بایگانی سیستم فایل محلی. فایل اصلی مورد خزش قرار می‌گیرد و صفحه‌ها و سایر فایل‌ها به صورت انفرادی روی سیستم فایل بایگانی ذخیره می‌شوند. دستیابی بوسیله شناوری مستقیم در نظام فایل انجام می‌شود.

استفاده ترجیحی

این روش برای بایگانی سایت شرکت ها یا سازمانی و بایگانی که در مقیاس کوچکی توسعه یافته، توصیه شده است. بسته به استفاده از این بایگانی، صحت مبحث باید به دقت مورد بررسی قرار گیرد مخصوصاً برای بایگانی سازمانی برای بایگانی با مقیاس کوچک، تعادل بین سختی برای سازماندهی مداوم مجموع فایل ها و سادگی دسترسی ایجاد شده توسط این رویکرد باید بر اساس «مورد به مورد» ارزیابی گردد. برای بایگان های وب با مقیاس بزرگ و متوسط از این روش باید اجتناب شود.

ابزارها

این راهبرد برای بایگانی وب در مقیاس کوچک و متوسط ساده ترین رویکرد است که با بسیاری از ابزار های موجود مانند HTTrack قابل دسترس است.

2-3-4-بایگانی های مبتنی بر خدمت وب

اشاره

*بایگانی های مبتنی بر خدمت وب (1)

با وجود تقاضاهای بیشتر، این گزینه تطبیق بهتری را برای نام گذاری منبع و ساختار اسناد ایجاد می کند (شکل 6-1). همچنین، برای اجتناب از محدودیت های اندازه نظام فایل که برای بایگانی وب در مقیاس بزرگ بحرانی است، مجاز می کند.

عکس

استفاده ترجیحی

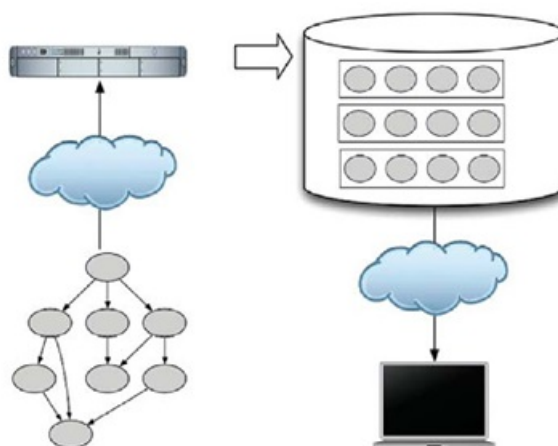
این روش برای بایگانی سایت شرکت ها یا سازمانی و بایگانی که در مقیاس کوچکی توسعه یافته، توصیه شده است. بسته به استفاده از این بایگانی، صحت مبحث باید به دقت مورد بررسی قرار گیرد مخصوصاً برای بایگانی سازمانی. برای بایگانی با مقیاس کوچک، تعادل بین سختی برای سازماندهی مداوم مجموع فایل ها و سادگی دسترسی ایجاد شده توسط این رویکرد باید براساس «مورد به مورد» ارزیابی گردد. برای بایگان‌های وب با مقیاس بزرگ و متوسط از این روش باید اجتناب شود.

ابزارها

این راهبرد برای بایگانی وب در مقیاس کوچک و متوسط ساده ترین رویکرد است که با بسیاری از ابزار های موجود مانند HTTrack قابل دسترس است.

۴-۳-۲-بایگانی‌های مبتنی بر خدمت وب^۱

با وجود تقاضاهای بیشتر، این گزینه، تطبیق بهتری را برای نام گذاری منبع و ساختار اسناد ایجاد می کند (شکل ۶-۱). همچنین، برای اجتناب از محدودیت‌های اندازه نظام فایل که برای بایگانی وب در مقیاس بزرگ بحرانی است، مجاز می کند.



شکل ۶-۱. مدل مبتنی بر خدمت وب. سایت اصلی مورد خزش قرار می گیرد و پاسخ هادر مخزن بدون تغییر ذخیره می شوند (فایل های) که اجتناب از پیمان نامه مسیریابی نام گذاری فایل های نظام و تغییر ساختار پیوند را اجازه می دهند. دستیابی نیازمند یک سرور وب است که محتوا را در مخازن واگشی کند و آن را به عنوان یک پاسخ به کاربر بفرستد.

1. Web-served archive

شکل ۶-۱. مدل مبتنی بر خدمت وب. سایت اصلی مورد خزش قرار می گیرد و پاسخ ها در مخزن بدون تغییر ذخیره می شوند (فایل های) که اجتناب از پیمان نامه مسیریابی نام گذاری فایل های نظام و تغییر ساختار پیوند را اجازه می دهند. دست یابی نیازمند یک سرور وب است که محتوا را در مخازن واگشی کند و آن را به عنوان یک پاسخ به کاربر بفرستد.

این روش، بر اساس بایگانی پاسخ است (در مقایسه با اولی که بر اساس بایگانی فایل است). پاسخ ها از سرور اصلی (منبع) بدون تغییر در فایل های مخزن (WARC 1) ذخیره می شوند که اجازه می دهد تا پشت سر آخرین تا کاربر بایگانی با یک سرور HTTP خدمت رسانی را انجام دهد.

پیشینه های یک فایل WARC (فایل آرشیوی مبتنی بر خدمت) رشته ای از فایل های وب را ثبت می کند هر صفحه به وسیله یک (برچسب) سرآمد که به طور مختصر محتوای حاصل شده و طولش توصیف کند، پیشی می گیرد. به علاوه، محتوای اولیه ثبت شده و WARC (فایل آرشیوی مبتنی بر خدمت) محتوای ثانویه را نیز در بر می گیرد مانند ابر داده و نقل و انتقالات فایل اصلی (منبع). اندازه یک فایل WARC (فایل آرشیوی مبتنی بر خدمت) می تواند تا حدود صدها مگابایت فرق داشته باشد. هر پیشینه یک نقطه شروع دارد که دسترسی مستقیم به پیشینه های انفرادی را بدون بارگیری و تجزیه تمام فایل های WARC (فایل آرشیوی مبتنی بر خدمت) انجام می دهد. نقاط شروع پیشینه های انفرادی در یک نمایه مرتب شده توسط URL ذخیره می شوند. از این رو، به سرعت برای پیشینه های انفرادی که بر اساس URL شان خارج از یک مجموع فایل WARC (فایل آرشیوی مبتنی بر خدمت) هستند، استخراج می شوند؛ که برای دسترسی در حالت شناوری تطبیق یافته اند. سپس ثبت هایی برای سرور وب تصویب می شوند که آن ها را سرویس گیرنده تهیه کرده است.

نگهداری و حفاظت از طرح نام گذاری شمای اصلی منبع (شامل پارامترهایی در صفحه های پویا)، شناوری در سایت را مجاز می کند همان طور که خزش شده است کاربر بایگانی می تواند از تمام سیرهای پیروی شده توسط خزشگر، بار دیگر پیمایش نماید (حرکت کند).

توضیح

مزیت عمده استفاده از مخزن های WARC امکانات غلبه یافتن بر محدودیت سیستم فایل ذخیره سازی در رابطه با اندازه (تعداد کمی از فایل های انفرادی در نهایت در نظام فایل بایگانی ذخیره می شوند) و فضای نام (نام گذاری انفرادی فایل های وب می تواند حفظ و نگهداری شود) می باشد. دستیابی به بایگانی اینترنت از طریق Wayback Machin (که دسترسی به 500tb از مجموعه وب را می دهد) نشان می دهد که این رویکرد به نسبت ثابت افزایش داشته است نه مانند دیگران. اشکال این رویکرد این است که دسترسی مستقیم به فایل های ذخیره شده غیر ممکن است. دو لایه اضافی استفاده شده برای دسترسی به محتوا ضروری هستند:

نظام نمایه فایل WARC و سرور وب. این دو لایه، پیچیدگی برجسته ای ندارند، اما برای دستیابی به محیط به اجرا نیاز دارد که می تواند برای راه اندازی و نگهداری در سازمان های کوچک سخت باشد. این میانجی گر می تواند همچنین مشکلاتی را برای انتقال محتوا افزایش دهد چون نیاز دارد که ساز و کارهای

پیوند به طور مناسب از محیط وب زنده به محیط بایگانی، مسیره‌ی نمایندگی (فرض می‌کنیم که پیوندهای اصلی بدون تغییر در بایگانی نگهداری می‌شوند که این مزیت مهم این روش است). این کار می‌تواند در سطح نمایش صفحه و در سطح نماینده انجام شود.

گزینه نخست در افزودن به صفحه فرستاده شده به مرورگر کاربر بایگانی را در بر می‌گیرد که متن سند این کار را خواهد کرد، در خزش تفسیر مجدد پیوندها در صفحه ای به نقطه ای در بایگانی انجام خواهد شد (یا آن‌ها را در یک شکل مرتبط تفسیر می‌کند) بایگانی اینترنت برای مثال این کار را با پیروی از کد جاوا-اسکریپت (1) هر صفحه فرستاده شده برای کاربران پیوست شده انجام می‌دهد.

```
<>>SCRIPT language="Javascript<
```

```
<!--
```

```
FILE ARCHIVED ON 20050308085053 AND RETRIEVED FROM THE//
```

```
.INTERNET ARCHIVE ON 20060514055212//
```

```
JAVASCRIPT APPENDED BY WAYBACK MACHINE, COPYRIGHT//
```

```
.INTERNET ARCHIVE
```

```
ALL OTHER CONTENT MAY ALSO BE PROTECTED BY COPYRIGHT//
```

```
.U.S.C(17
```

```
)).SECTION 108(a)(3//
```

```
!";var sWayBackCGI ="http://web.archive.org/web/20050308085053
```

```
)function xLateUrl(aCollection, sProp
```

```
var i = 0( )
```

```
++)for(i = 0; i < aCollection.length; i
```

```
if (aCollection[i][sProp].indexOf("mailto:") == -1
```

```
)aCollection[i][sProp].indexOf("javascript:") == -1
```

```
];aCollection[i][sProp] = sWayBackCGI + aCollection[i][sProp
```

```
var i = 0 }
```

```
++)for(i = 0; i < aCollection.length; i  
  
if (aCollection[i][sProp].indexOf('mailto:') == -1  
  
)aCollection[i][sProp].indexOf('javascript:') == -1  
  
];aCollection[i][sProp] = sWayBackCGI + aCollection[i][sProp  
  
var i = 0 }
```

```
++)for(i = 0; i < aCollection.length; i
```

```
if (aCollection[i][sProp].indexOf('mailto:') == -1
```

ص: 39

```
)aCollection[i][sProp].indexOf('javascript:') == -1
```

```
];aCollection[i][sProp] = sWayBackCGI + aCollection[i][sProp
```

```
}
```

```
");if (document.links) xLateUrl(document.links, "href
```

Web Archiving: Issues and Methods 35 1

```
");if (document.images) xLateUr (document.images, "src
```

```
");if (document.embeds) xLateUrl(document.embeds, "src
```

```
+if (document.body document.body.background( document.body.background = sWayBackCGI
```

سند. بدنه. زمینه

```
//->
```

```
>SCRIPT</
```

```
>HTML</
```

مشکل این روش این است که برخی پیوندها (جاسازی شده در متن سند) [\(1\)](#) تفسیر نخواهد شد و از این رو، در نقطه ای در وبگاه منبع می مانند. در برخی موارد، تفسیر کد صفحه، برخی رفتارها را فعال می کند، مانند تغییر مسیر حتی قبل از اینکه کد پیوست شده به عنوان مرورگر جدید تفسیر گردد. منتظر نباشید تا سند کامل برای تفسیر به دست آید و آن را نمایش دهد.

با استفاده از نماینده ای که تمام درخواست ها از مرورگر کاربر به بایگانی تغییر مسیر می دهد، بسیار مؤثرتر و کارآمد می باشد همانطور که مسیر دهی بعد از تفسیر پیوند رخ می دهد توسط تعامل کاربر (با کلیک کردن) و مرورگر انجام می شود و کد را برای تولید درخواست مناسب تفسیر می کند (HTML سند متن سرویس گیرنده جانبی فرمت های دیگر) این کارآمدترین به عنوان ظرفیت مرورگر اصلی برای تفسیر مجموعه کدهای استاندارد برای کاری که معمولاً بر روی وب انجام می شود. این رویکرد، نیاز به ایجاد یک نماینده دارد که به بایگانی تغییر مسیر دهد و به یک پارامتر که یک مرورگر از آن استفاده کند تا بتواند تقاضاهای زیادی برای یک محیط بایگانی پیوسته داشته باشد. استفاده از اتصال مرورگر برای مدیریت انتقال شکل باز به محیط نماینده می تواند این کار را برای کاربران نهایی آسان سازد.

استفادهٔ ترجیحی

این روش برای بایگانی مقیاس بزرگ و متوسط و همچنین بایگانی های کوچک تر مناسب است که در زمینهٔ حفاظت از صحت محتوا هستند. چون این روش ها، پاسخ ها را از سرور منبع همانطور که از سرویس گیرنده دریافت شده است ذخیره می کنند بدون تغییر و

1- ساختار الگو گونه برای نمایش ترتیب حوادث: Script

بیشتری را از روش های دیگر فراهم می کند. چون بر سازماندهی فایل محلی بستگی دارد، همچنین مناسب برای انتقال و همچنین تحویل محتواست.

ابزارها

این روش، به یک زیر ساخت دسترسی و به یک خزنده بایگانی (مانند HERITRIX) و یک نظام نمایه سازی برای فایل های WARC نیاز دارد. IIPC زنجیره کاملی از ابزارهایی دارد که برای ایجاد این توابع عملیاتی توسعه یافته است.

3-3-4-بایگانی غیروب

توصیف

در این رویکرد که در شکل 1-7 ترسیم شده است، اسنادی که بر روی وب بوده و از متن ابر متن استخراج شده و در یک سبک متفاوت برای دسترسی منطقی و /یا فرمت سازمان دهی مجدد می شوند.

این می تواند موردی باشد که وقتی یک مجموعه اسناد از وب برداشته می شوند، از دسترسی منطقی مبتنی بر پیوند به یک پیوند مبتنی بر فهرست دوباره سازماندهی می شوند.

عکس

بیشتری را از روش‌های دیگر فراهم می‌کند. چون بر سازماندهی فایل محلی بستگی دارد، همچنین مناسب برای انتقال و همچنین تحویل محتواست.

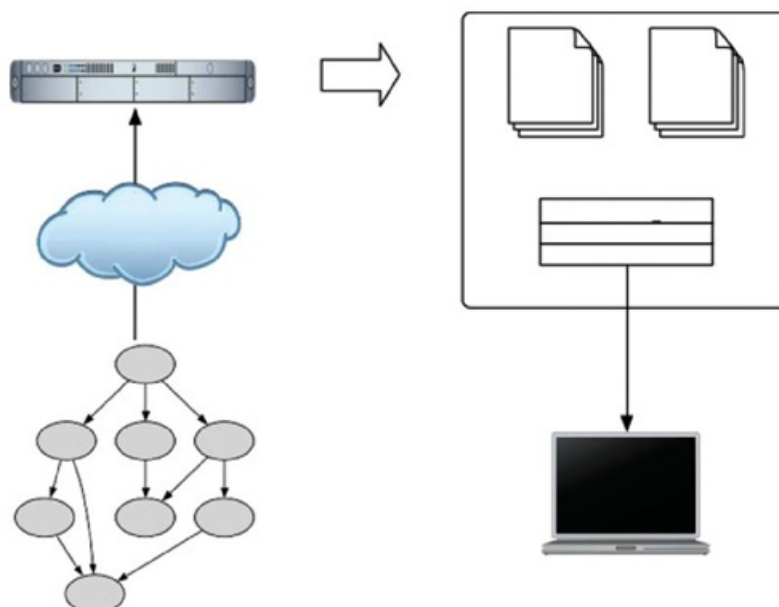
ابزارها

این روش، به یک زیر ساخت دسترسی و به یک خزنده بایگانی (مانند HERITRIX) و یک نظام نمایه سازی برای فایل‌های WARC نیاز دارد. IIPC زنجیره کاملی از ابزارهایی دارد که برای ایجاد این توابع عملیاتی توسعه یافته است.

۴-۳-۳- بایگانی غیرواب

توصیف

در این رویکرد که در شکل ۷-۱ ترسیم شده است، اسنادی که بر روی وب بوده و از متن ابر متن استخراج شده و در یک سبک متفاوت برای دسترسی منطقی و/یا فرمت سازماندهی مجدد می‌شوند. این می‌تواند موردی باشد که وقتی یک مجموعه اسناد از وب برداشته می‌شوند، از دسترسی منطقی مبتنی بر پیوند به یک پیوند مبتنی بر فهرست دوباره سازماندهی می‌شوند.



شکل ۷-۱. اسناد از سایت اصلی در بایگانی سازماندهی مجدد می‌شوند، ساختار غیر وبی را پیروی می‌کنند، به طور موردی از فهرستی که دسترسی به اسناد انفرادی را فراهم می‌کند، استفاده می‌کند.

شکل ۷-۱. اسناد از سایت اصلی در بایگانی سازماندهی مجدد می‌شوند، ساختار غیر وبی را پیروی می‌کنند، به طور موردی از فهرستی که دسترسی به اسناد انفرادی را فراهم می‌کند، استفاده می‌کند.

همچنین، این موردی است که وقتی یک صفحه یا حتی یک وبگاه کامل به فرمت PDF تغییر شکل یافته است. (1 Adobe Acrobat) این توابع عملیاتی را دارد و می تواند تمام یک وبگاه را در یک سند واحد PDF تغییر شکل دهد. در این مورد سند به طور مجازی چاپ شده است که یک انتقال بی حرکت و یک صفحه یاد داشت مانند سازماندهی را شامل می شود، حتی اگر پیوند دهی بتواند هنوز با استفاده از شمای داخلی نام گذاری مناسب کار کند.

توضیح

این رویکرد اساساً برای دریافت موضوع هایی که در اصل ایجاد شده اند، به طور مستقل از وب سازماندهی شده اند، ایجاد شده است. این مورد برای مجموعه های بزرگی از کتاب های رقومی به عنوان نمونه، مقاله ها، موزیک، و ویدئوهای ساخته شده که بر روی وب موجود هستند، ولی سازماندهی اصلی شان ابرمتنی نمی باشد چون مبتنی بر کاتالوگ هستند.

این رویکرد می تواند در این مورد برای چسبیدن به معماری اطلاعات اصلی و بایگانی این مجموعه ها با یکدیگر با کاتالوگ های ادغام شده در کاتالوگ بایگانی، مرجع باشد. فرض شده است که محتوای ابر متن بدون ربط پنداشته شده است و می تواند منتشر شود.

برای مثال این مورد در طرح انبارهای الکترونیکی kb در نیوزیلند جایی که انتشارات علمی الزویر (2) اجرا شده است که در یک نظام مبتنی بر فهرست بایگانی شده است. حقیقت این است که الزویر به این مطالب دسترسی داشته است و به عنوان یک مجموعه کمکی برای محتوای خودش مورد بررسی قرار گرفته و به عنوان انتشار علمی سنتی ساخت یافته است.

مرجع

اشاره

این روش برای مجموعه های محتوا و نه ساخت یافته در حالت وب دلالت می کند.

4-3-4 خلاصه

جدول 2-1- انواع گوناگونی از بایگانی های وب، موارد مرجع های مورد استفاده، ابزارهای و مزیت های و عدم مزیت ها را به طور خلاصه شرح داده است.

ص: 42

1- نرم افزاری از محصول شرکت Adobe برای ساخت و تهیه فایل های (PDF)

Elsevier-2

بایگانی شبکه وب: مباحث و روش ها ۴۳

جدول ۱-۲. خلاصه ای از انواع بایگانی وب

نوع بایگانی	نظام فایل محلی	مبتنی بر خدمت	غیروپ
توصیف	تمام پیوند به پیوندهای مرتبط (رابطه) تبدیل شده‌اند. و شناوری ابر متن، به‌طور قطع بر روی نظام‌فایل محلی انجام شده است.	سرور یک وب، برای دسترسی از طریق هر سندی که به‌کار رفته است نصب شده است و شناوری ابر متن به بایگانی اصلی بسته شده است.	اسناد از محتوای ابر متن اصلی استخراج شده و مجدداً برای یک منطق مختلف سازماندهی شده‌اند.
استفاده مرجع	بایگانی سایت منفرد و بایگانی در مقیاس متوسط و کوچک	بایگانی در مقیاس کوچک و متوسط	بایگانی مجموعه‌های خاص (غیروپ)
ابزارها	کپی کننده وبگاه (مانند HTTrack)	خزنده بایگانی (مانند Heritrix) و نظام نمایه‌سازی برای فایل‌های WARC	بستگی به ساختار محتوای نهایی دارد.
مزیت‌ها	برای اجرا ساده است	صحت، قابلیت مقیاس پذیری	قادر به ایجاد یکپارچگی در فهرست‌های سنتی یا سایر معماری‌های اطلاعاتی محلی
معایب	افزایش مقیاس ندارد. نیاز به نامگذاری مجدد و محدود شدن سازماندهی مجدد محتوا برای شناوری ابر متن‌ها است. نیاز به مدیریت در سطح نظام فایل در مجموعه بایگانی و نگارش موارد دارد.	اجرا در غیاب نرم افزار یکپارچه‌سازی سخت می‌باشد (این کار ممکن است در آینده تغییر کند).	فقدان ساختار ابر متن. فقط می‌تواند برای اسناد مجزا و غیروپ به‌کار برده شود.

۴-۴- کیفیت و تمامیت (کامل بودن)

به‌طور کلی، کیفیت می‌تواند در یک حالت وظیفه‌ای (متناسب برای استفاده خاص) یا در یک حالت هدفمند (منطبق با ویژگی‌های اندازه‌گیری) تعریف شود. اصطلاح کیفیت، برای مجموعه فرهنگی در مضمون‌های مختلف به‌کار برده شده است. شخص می‌تواند از آن برای کنترل وضعیت حفاظت، کامل نمودن موارد یا مجموعه‌ها، سطح محتواهای هوشمند و علمی و غیره استفاده شود. در هر مورد، آن با مقیاس مطلوب کامل در یک ناحیه خاص و (حفاظت فیزیکی، پوشش یک قلمرو، صحت انتخاب) مرتبط می‌باشد.

برای بایگانی‌های وب همان‌طور که دیده شده، بیشترین نقایص ناشی از سختی برای جمع‌آوری محتوا از طریق واسط HTTP (بخش قبلی در مورد بایگانی سرویس گیرنده جانبی را ببینید) و سختی تحویل در یک حالت منسجم از نتایج محتوا (بخش «سازماندهی و ذخیره‌سازی» را ببینید) می‌باشد. از این رو، کیفیت بایگانی وب در این فصل به‌عنوان ۱) تکامل مطالب (فایل‌های پیوندی) بایگانی شده درون فضای پیرامون هدف، و ۲) توانایی برای تحویل شکل اصلی سایت، به‌ویژه در رابطه با هدایت و تعامل با کاربران (میزانس، ۲۰۰۵) مورد بررسی قرار گرفته است.

جدول ۱-۲. خلاصه ای از انواع بایگانی وب

۴-۴- کیفیت و تمامیت (کامل بودن)

به طور کلی، کیفیت می تواند در یک حالت وظیفه ای (متناسب برای استفاده خاص) یا در یک حالت هدفمند (منطبق با ویژگی های اندازه گیری) تعریف شود. اصطلاح کیفیت، برای مجموعه فرهنگی در مضمون های مختلف به کار برده شده است. شخص می تواند از آن برای کنترل وضعیت حفاظت، کامل نمودن موارد یا مجموعه ها، سطح محتواهای هوشمند و علمی و غیره استفاده شود. در هر مورد، آن با مقیاس مطلوب کامل در یک ناحیه خاص و (حفاظت فیزیکی پوشش یک قلمرو، صحت انتخاب) مرتبط می باشد.

برای بایگانی های وب همانطور که دیده شده، بیشترین نقایص ناشی از سختی برای جمع آوری محتوا از طریق واسط HTTP (بخش قبلی در مورد بایگانی سرویس گیرنده جانبی را ببینید) و سختی تحویل در یک حالت منسجم از نتایج محتوا (بخش «سازماندهی و ذخیره سازی» را ببینید) می باشد.

از این رو، کیفیت بایگانی وب در این فصل به عنوان 1) تکامل مطالب (فایل های پیوندی) بایگانی شده درون فضای پیرامون هدف و 2) توانایی برای تحویل شکل اصلی سایت به ویژه در رابطه با هدایت و تعامل با کاربران (میزانس، 2005) مورد بررسی قرار گرفته است.

تکامل می تواند به طور افقی با تعدادی از نقاط ورودی مربوط که در فضای پیرامون طراحی شده، با نمایش هندسی ارزیابی شود و از نظر برنامه کاربردی عمودی با تعدادی از گره های پیوندی می توان که از این نقطه ورودی دریافت شده اند اندازه گیری می شود. معمولا نقاط ورودی صفحه های اصلی سایت هستند و پیوندها می توانند کاربران را به هر یک از نقاط ورودی هدایت کنند (سایت دیگر) یا به عناصر همان سایت بفرستند. این مورد برای بایگانی مبتنی بر سایت است.

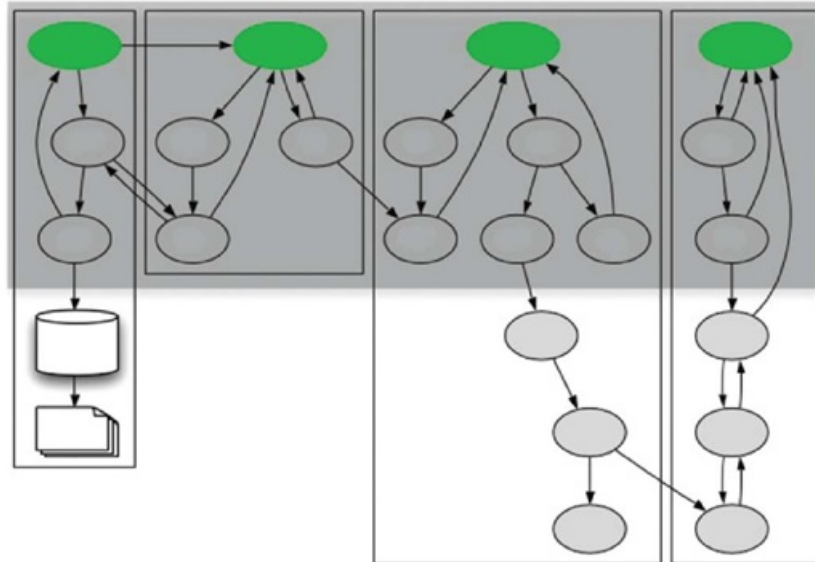
به هر حال، در برخی موارد، گنجایش برنامه کاربردی عمودی برای جا سازی عناصر (برای مثال تصاویر همراه با یک صفحه) محدود شده است و مجموعه فقط به طور افقی با نادیده گرفتن سطح سایت، سازماندهی شده است. برای مثال این مورد برای خزش موضوع های خالص در جایی که صفحه ها را در بر نمی گیرد بر اساس تعلق شان به سایت می باشد، اما فقط بر روی ارتباط شان با موضوع است.

به طور مطلوب، بایگانی وب باید به طور عمودی همچنین افقی کامل شود. اما در عمل، برای دسترسی سخت است و اولویت هایی باید قرار داده شوند. بایگانی زمانی که کامل شدن افقی به کامل شدن عمودی ترجیح داده شود، «بسیط» نامیده می شود.

عکس

تکامل می‌تواند به‌طور افقی با تعدادی از نقاط ورودی مربوط که در فضای پیرامون طراحی شده، با نمایش هندسی، ارزیابی شود و از نظر برنامه کاربردی عمودی، با تعدادی از گره‌های پیوندی می‌توان که از این نقطه ورودی دریافت شده‌اند، اندازه‌گیری می‌شود. معمولاً، نقاط ورودی، صفحه‌های اصلی سایت هستند و پیوندها می‌توانند کاربران را به هریک از نقاط ورودی هدایت کنند (سایت دیگر) یا به عناصر همان سایت بفرستند. این مورد برای بایگانی مبتنی بر سایت است.

بمهر حال، در برخی موارد، گنجایش برنامه کاربردی عمودی برای جا سازی عناصر (برای مثال تصاویر همراه با یک صفحه) محدود شده است و مجموعه فقط به‌طور افقی با نادیده گرفتن سطح سایت، سازماندهی شده است. برای مثال، این مورد برای خزش موضوع‌های خالص در جایی که صفحه‌ها را در بر نمی‌گیرد بر اساس تعلقشان به سایت می‌باشد، اما فقط بر روی ارتباطشان با موضوع است. به‌طور مطلوب، بایگانی وب باید به‌طور عمودی همچنین افقی کامل شود. اما در عمل، برای دسترسی سخت است و اولویت‌هایی باید قرار داده شوند. بایگانی زمانی که کامل شدن افقی به کامل شدن عمودی ترجیح داده شود، «بسیط» نامیده می‌شود.



شکل 1-8 مجموعه‌های بسیط شامل بیشتر سایت‌هاست ولی فقط در سطحی بایگانی می‌شود. فقط محتوا در ناحیه سایه‌دار بایگانی می‌شود. صفحه‌ها عمیق در سلسله مراتب (سی 6، سی 7، سی 8، دی 3، دی 3، دی 5) و نیز محتوایی که در پایگاه داده پنهان شده (وب پنهان) تصرف خواهند شد.

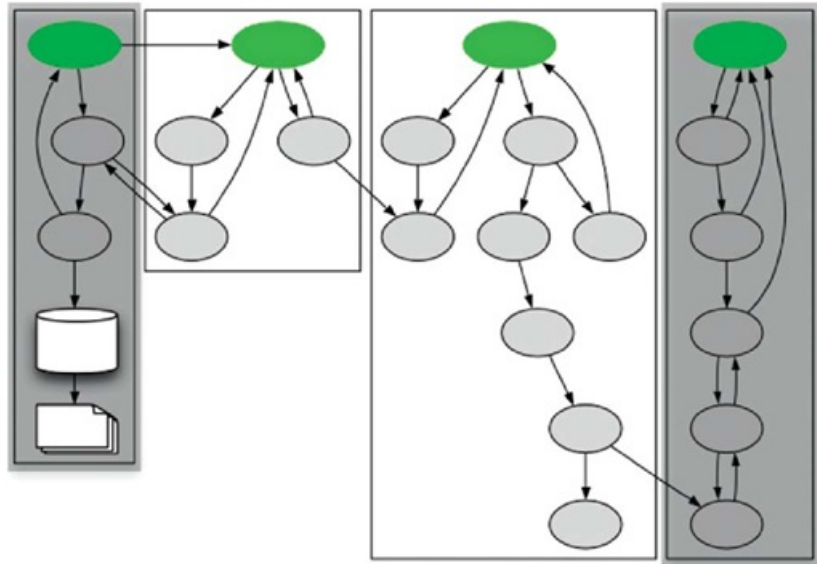
شکل 1-8. مجموعه‌های بسیط شامل بیشتر سایت‌هاست ولی فقط در سطحی بایگانی می‌شود. فقط محتوا در ناحیه سایه دار بایگانی می‌شود. صفحه‌ها عمیق در سلسله مراتب (سی 6، سی 7، سی 8، دی 3، دی 3، دی 5) و نیز محتوایی که پایگاه داده پنهان شده (وب پنهان) تصرف خواهند شد.

برای مثال، این مورد برای مجموعه بایگانی اینترنت است که به وسیله الکسا (همچنین برنر، 1997؛ کیمپتون و همکاران، 2006) ارائه شده است، خزشگر الکسا یک رویکرد خزش سطح اول را استفاده می کند و عمق خزش را برای یک سایت بر طبق ترافیک اندازه گیری شده برای این سایت وفق می دهد.

بر عکس، بایگانی «متمرکز» نامیده می شود زمانی است که تکامل عمودی به کامل شدن افقی ترجیح داده می شود (شکل 9-1 را ببینید).

عکس

برای مثال، این مورد برای مجموعه بایگانی اینترنت است که به وسیله الکسا (همچنین برنر، ۱۹۹۷؛ کیمپتون و همکاران، ۲۰۰۶) ارائه شده است، خزشگر الکسا یک رویکرد خزش سطح اول را استفاده می‌کند و عمق خزش را برای یک سایت بر طبق ترافیک اندازه‌گیری شده برای این سایت وفق می‌دهد. برعکس، بایگانی «متمرکز» نامیده می‌شود زمانی است که تکامل عمودی به کامل شدن افقی ترجیح داده می‌شود (شکل ۹-۱ را ببینید).



شکل ۹-۱. بایگانی بسیط سایت‌های کمتری مورد خزش قرار می‌گیرند ولی خزش در عمق انجام می‌شود. ولی سایت الف و دی مورد بایگانی قرار می‌گیرند ولی در حالت کلی شامل بخش وب پنهان سایت الف.

برای مثال، در این مورد، وقتی اولویت نخست - سایت برای خزش گرها استفاده شده یا وقتی که یک بازبینی دستی، در جایی که مورد نیاز است، بر روی آن انجام می‌شود، بایگانی تکمیلی انجام می‌شود. بایگانی متمرکز حتی متقاضی بیشتری برای سایت‌های وب پنهان دارد (همچنین، سایت‌های وب عمیق نامیده شده است) جایی که دسترسی به تمام محتوا با خزش گرها ممکن نیست.

۵ - بررسی عمومی مراحل اولیه جاری

بایگانی وب می‌تواند در چندین روش طبقه‌بندی شود. در این بخش، روش‌های مهم را بررسی خواهیم کرد و این فرصت را برای ارائه برخی مراحل اولیه بایگانی وب و مقایسه با رویکردهای گوناگون نشان می‌دهیم.

شکل 1-9. بایگانی بسیط سایت‌های کمتری مورد خزش قرار می‌گیرند ولی خزش در عمق انجام می‌شود. ولی سایت الف و دی مورد بایگانی قرار می‌گیرند ولی در حالت کلی شامل بخش وب پنهان سایت الف.

برای مثال، در این مورد، وقتی اولویت نخست - سایت برای خزش گرها استفاده شده یا وقتی که یک بازبینی دستی، در جایی که مورد نیاز است، بر روی آن انجام می‌شود، بایگانی تکمیلی انجام می‌شود. بایگانی متمرکز حتی متقاضی بیشتری برای سایت‌های وب پنهان دارد (همچنین، سایت‌های وب عمیق نامیده شده است) جایی که دسترسی به تمام محتوا با خزش گرها ممکن نیست.

5- بررسی عمومی مراحل اولیه جاری

بایگانی وب می تواند در چندین روش طبقه بندی شود. در این بخش روش های مهم را بررسی خواهیم کرد و این فرصت را برای ارائه برخی مراحل اولیه بایگانی وب و مقایسه با رویکردهای گوناگون نشان می دهیم.

ص: 45

نوع ایجاد سازماندهی و میزبانی از بایگانی، نخستین معیار برای طبقه بندی بایگانی های وب می باشد. برخی دسترسی عمومی به مجموعه های شان را فراهم می کنند (WA عمومی) برخی این کار را نمی کنند (WA خصوصی (یا مخفی)).

در میان بایگانی های عمومی وب، برخی دسترسی پیوسته را فراهم می کنند برخی دسترسی به سایت را در اتاق های مطالعه فراهم می کنند (بایگانی وب عمومی پیوسته و بایگانی وب عمومی غیر پیوسته). همچنین برخی و در بیشتر موارد در ابتدا مجموعه غیر رقومی را مدیریت می کنند (بایگانی وب هیبرید) و بالاخره برخی در حالت سرمایه گذاری یا بدون منفعت (بایگانی وب غیر تجاری) با در نظر گرفتن اینکه برخی شرکت های تجاری هستند (بایگانی وب تجاری).

مؤسسه های میراث سنتی (کتابخانه ها، بایگانی ها، موزه ها) که مجموعه هایشان را برای وب گسترش داده اند، با یکدیگر بیشترین قسمت مقوله بایگانی وب هیبرید عمومی را تشکیل داده اند. کتابخانه های ملی چندین کشور متعلق به این مقوله هستند (سوئد و استرالیا که نخستین بار به 1996 بر می گردد و اکنون کشورهای زیادی هستند) (1).

بایگانی های ملی، منطقه ای و شهری همچنین شروع به بایگانی وبگاه های مجاز دولتی و محلی نموده اند (2). سازمانی بر روی شکل های جدید صنعتی مانند v2 مستقر در روتردام هلند کار می کند که در حال یکپارچه سازی شبکه در یک انعکاس و عملکردی عمومی برای حفاظت از رسانه ناپایدار می باشد (فوکونیر و فرومه 2004) (3).

تمام این بایگانی ها می توانند به عنوان آرشیو وب عمومی هیبرید رده بندی شوند. همان طور که محتوای وب را در یک مضمون بزرگ تر از مجموعه ها مجتمع می کنند. بیشتر آن ها فقط دسترسی غیر پیوسته برای اسناد را در حال حاضر فراهم می کنند.

از میان آن ها، کتابخانه اسکندریه در مصر یکی از معدود دسترسی های پیوسته برای مجموعه بایگانی وب است (معکوس نمودن بایگانی اینترنتی) و یک مثال از بایگانی وب عمومی هیبرید غیر تجاری پیوسته می باشد.

نفوذپذیری اینترنت همچنین ضرورت برخی از انواع جدید سازمان های بایگانی را مجاز می داند که فقط مجموعه رقومی و تهیه دسترسی پیوسته را نگه می دارد که به عنوان بایگانی وب پیوسته غیر تجاری عمومی، رده بندی خواهند شد.

بایگانی اینترنت، مثال مهمی در این مقوله می باشد (فصل 9 را ببینید) (کیمپتون و همکارانش، 2006).

ص: 46

1- با پیروی از آن ها، چندین کتابخانه ملی بایگانی وب را شروع کرده و برنامه هایی را اجرا می کنند (این سیاهه جامع نیست): در اروپا: فنلاند، دانمارک، نروژ، ایسلند فرانسه، جمهوری چک، اسلونی، ایتالیا و یونان؛ در آسیا ژاپن چین و سنگاپور و کتابخانه کنگره در آمریکا با پیروی از آن ها، چندین کتابخانه ملی بایگانی وب را شروع کرده و برنامه هایی را اجرا می کنند (این سیاهه جامع نیست): در اروپا: فنلاند، دانمارک، نروژ، ایسلند فرانسه، جمهوری چک، اسلونی، ایتالیا و یونان؛ در آسیا ژاپن چین و سنگاپور و کتابخانه کنگره در آمریکا.

2- آرشیوهای ملی استرالیا (آرشیوهای ملی استرالیا، 2001) بریتانیا (براون، 2006)؛ کانادا ایالات متحده آمریکا (کارلین، 2006)

بایگانی وب نظام مند را شروع کرده اند. هم چنین طرح شهر Antwerp DAVID را ببینید (بودرس و آینده، 2002).

Fauconnier and Frommé -3

برخی شرکت های تجاری مجموعه های بزرگی از محتوای وب عمومی را بایگانی می کنند مانند گوگل با نهانگاهش (1) و بایگانی هانزو (2). مثال هایی از بایگانی وب تجاری عمومی پیوسته هستند.

سرانجام، بسیاری از سازمان ها، بایگانی وب داخلی را برای اهداف شخصی توسعه می دهند که به عنوان بایگانی وب خصوصی (مخفی) طبقه بندی خواهند شد. منحصر به فرد بودن نوع دسترسی (پیوسته یا غیر پیوسته) و همچنین وضعیت تجاری در اینجا ارتباط کمتری دارد چون این بایگانی ها فقط برای استفاده خصوصی هستند.

2-5-2- حوزه (دامنه)

روش مفید دیگر برای دسته بندی بایگانی های وب، بررسی حوزه ای است که آن ها اتخاذ می کنند. بایگانی های وب می توانند هر یک، سایت، سر فصل ها یا متمرکز بر حوزه باشند.

5-2-1- بایگانی مرکزی سایت

این نوع بایگانی، بر روی یک سایت مشخص متمرکز شده است که تقریباً به وسیله و برای ایجاد کننده سایت اجرا شده است. این حوزه بندی، از این رو، تقریباً برای بایگانی وب خصوصی استفاده شده است. مثلاً بسیاری از شرکت ها، مسئول تمام محتواهایی هستند که منتشر می کنند و باید مطمئن باشند که می توانند به روش های قدیمی تر سایت و بنوشت ها مراجعه کنند. این نوع بایگانی ترجیحاً از کپی کنندگان سایت و برخی تهیه کنندگان خدمات اینترنتی استفاده می کند که برای این نوع بایگانی داخلی مناسب، به وجود آمده اند (3).

5-2-2- بایگانی مرکزی عنوان

بایگانی های وب، عمومی و عمومی تر شده اند، اغلب به وسیله نیازهای پژوهشی مستقیم اجرا می شوند تا موقعی که کار بر روی یک فیلد مشخص و انعکاس بر روی وب انجام می شود، دانش پژوهان زیادی با طبیعت بی دوام انتشار وب روبه رو می شوند، جایی که طول عمر وبگاه برای بازبینی علمی (تکذیب نیازمند دستیابی به همان داده است) و همچنین برای ارجاع با دوام، نامناسب می باشد.

به این دلیل است که طرح های مختلف اغلب در کتابخانه های دانشگاه میزبانی می شوند و تحت حفاظت مطالب اولیه برای پژوهش قرار می گیرند، مانند بایگانی رومی برای مطالعات چینی در دانشگاه هایدلبرگ (4) در آلمان یا آرچیپل (5) برای تحلیل سایت های سیاسی هلندی در دانشگاه گرونینگن (6) در هلند

ص: 47

Cashe -1

Hanzo -2

3- برای نمونه hanzoarchives.com را ببینید.

Digital Archive for Chinese Studies (DACHS) at Heidelberg University -4

Archipol -5

Groningen -6

(ورمن 1) و همکاران، 2002) این پروژه ها نه تنها جهت دهی یک عنوان را به اشتراک می گذارند، بلکه از

یک شبکه آگاهی دهنده نیز استفاده می کنند.

از طرفی دیگر، پژوهشگران که تغذیه های صحیح و به روز

طرح های متمرکز دیگری در کتابخانه ها از طریق جست و جوی فعال و بایگانی وبگاه گزینشی (2) انجام شده است و مانند طرح مینروا (3) از کتابخانه کنگره (اشنایدر و همکاران 2003) یا بایگانی وب انتخابات فرانسه که توسط کتابشناسی ملی فرانسه ایجاد شده است (میزانس، 2005). در مقایسه با رویکرد قبلی مبتنی بر عنوان متمرکز، کشف سایت ها به طور طبیعی به عنوان محصول فرعی از فعالیت پژوهشی ایجاد نشده است. و به عنوان یک فعالیت خاص نیاز به بررسی دارد.

بالاخره، برخی طرح ها که با این مقوله مرتبط هستند. از خزش عنوان برای کشف و ضبط محتوا مرتبط با مطلب مشابه استفاده می کنند (چاکرابارتی (4) و همکاران، 1999؛ برگمارک 2002 (5)؛ برگمارک و همکاران، 2002؛ کیواین (6) و همکاران، 2004). کشف و فیلتر کردن خودکار با استفاده از فن سنتی خزش در ترکیب ارزیابی سطح صفحه در محتوای متنی انجام شده است که گاهی اوقات با کاویدن ساختارهای پیوندی آمیخته می شود. مجاورت با عنوان، می تواند از مجموعه ای از نوشتجات یا از باز خورد کاربر، فرا گرفته شود. اگر چه محتمل است، این ناحیه باز هم نیاز به پژوهش برای به کار برده شدن در بایگانی دارد.

3-2-5- بایگانی مرکزی حوزه ای

ساخت بایگانی می تواند بر اساس محل محتوا نیز انجام شود و به این ترتیب، نوع سوم در بایگانی وب حوزه ای را مشخص می کند. در اینجا کلمه «حوزه» (7) در جهت واژه شبکه یا به وسیله پسوند در جهت واژه ملی استفاده شده است که یک معیار ترکیبی برای سایت های هدف از یک کشور خاص می باشد (8).

نظام ملی حوزه ای یک گزینش ساده و قابل تعقیب قانونی در محتوا را بر اساس نام های دامنه اجازه می دهد. این حقیقتی است که نام های دامنه حتی برای سطوح بالایی حوزه ای توسط نمایندگی رسمی آی.سی.آن (9) مدیریت شده است، که واقعاً از قوانین در رابطه با نام دادن مشخصه عملیاتی و سازمان ها پیروی نمی کنند، ولی بیشتر سنت ها را پیروی می کنند (لیو و آلبیتر 1999) (10). همچنین می توانید در مورد تکامل بر روی نام های اینترنتی به (کوهلر، 1999) نگاه کنید. از این رو، شخص می تواند انواع پسوند های

ص: 48

Voerman -1

Electoral Web sites -2

Minerva -3

Chakrabarti -4

Bergmark -5

Qin -6

Domain -7

8- برای بحث در مورد روش ممکن مرزبندی کردن فضای اینترنت ملی آرویدسون و دیگران (2000)؛ ابایت بول و دیگران (2002)؛ لامپوس و دیگران (2004) را ببینید، برای مطالعاتی در مورد ویژگی های فضای اینترنت ملی بیزایتس و دیگران (2005 الف، 2005 ب) و گومس و سیلوا (2003) را ببینید.

ICANN -9

Liu and Albitz -10

عملیاتی یا عمومی مانند Com و edu و انواع سیستم اطلاعاتی جغرافیایی مانند (1) (ch.jp) و انواعی را در دامنه سطح نخست تشخیص دهد و که (اغلب دامنه سطح بالا نامیده می شود). دامنه های سطح بالای جغرافیایی اغلب بخش های فرعی عملیاتی دارند و (مانند asso.fr gov.mex). که به این معناست که دامنه سطح دوم (2) همچنین، در همان روش مدیریت خواهد شد. استثناها برای عملکردهایی مانند TLDها وجود دارد که بخش های فرعی جغرافیایی دیگر را دارا می باشند. (به وسیله ایالت ها). به هر حال، توجه داشته باشید که تمام این قسمت ها فضای حوزه اینترنتی است که توسط نمایندگی مدیریت شده اند (3). هر هویت تحت فرمان آن ها می تواند یک سیاست خاص را در رابطه با تخصیص و کنترل فضاهای شان به کار ببرد بنابراین رسیدن به بهره برداری از SLD, TLD برای گزینش بایگانی وب، در هر مورد طبق ارزیابی این سیاست بستگی دارد (مثلاً org.com.gTLD) به وسیله تمام انواع سازمان ها استفاده شده اند نه تنها توسط سازمان های تجاری برای com. و سازمان های بدون فایده برای org. چون محدودیتی برای ثبت نام وجود ندارد) به علاوه، برخی هویت ها تحت فرمان مدیریت TLD، سیاستشان را طبق زمان تغییر می دهند.

(org.net) با داشتن محدودیت ها قبل از سال 1996 استفاده شده اند و frTLD به طور قابل توجهی محدودیت ها را در سال 2005 کاهش داده است).

مزیت بزرگ دیگری که در اینجا باید ذکر کنیم، آوردن معیارهایی است که بتواند به طور خودکار توسط خزشگرها آشکار شوند مانند نام های دامنه پروژه های دیگری در واقع، رویکرد مرکزی دامنه را اجرا می نمایند. برخی، بر روی یک دامنه عمومی مانند (کروس و و همکاران، 2003؛ کارلین، 2004 و یا edu (لایل، 2004) متمرکز شده اند.

برخی از دامنه های ملی مانند ایالت کالتیورارو (4) در سال 1997 که توسط کتابخانه سلطنتی سوئد ایجاد شد، استفاده می کنند (آرویدسون 2000) که SeTLD و همچنین صفحه ها سوئدی پیوند یافته از آن و در مانند com. واقع شده است را پوشش می دهد.

3-5- روش های استفاده شده

طرح ها می توانند همچنین به طور قابل توجهی نسبت به رویکرد روش های بررسی که برای کشف اکتساب و توصیف محتوا در بر می گیرند و متفاوت باشند. یکی از تفاوت های مهم این است که محدوده ها در سراسر این مراحل، از دست به جای پردازش خودکار استفاده می کنند. اگر چه سادگی ظاهری این تضاد باید به عنوان پردازش خودکار متعادل می شود و در چندین سطح رخ می دهد (ضبط، استفاده از موتور جست و جو برای کشف دستی و غیره (میزانس، 2006 ب)، این مسئله، همچنین می تواند بایگانی وب را بر طبق این تضاد رده بندی نماید که به طور مستقیم اثر شدیدی بر روی قابلیت مقیاس پذیری و

ص: 49

1- ایزو 3166 دو حرف اول نام کشورها را اجازه می دهد به جز uk که باید gb باشد، و همچنین به جز این که اخیراً به سه حرف برای هر منطقه گسترش یافته است مثل حوزه cat در کاتولینا در اسپانیا.

SLD-2

3- در مورد حوزه دولتی و مفاهیم آن میولر (2002) را ببینید.

Kulturarw-4

همان طور که توانستیم پیش بینی کنیم، خود کار سازی این وظایف، در پایین آوردن فاحش هزینه های هر دستیابی به سایت توانا هستند (1). به طور مطلوب یک اپراتور تنها با اجرای خزش می تواند میلیون ها صفحه را کشف و بارگذاری کند. نظر به این که نمایه سازی تمام متن کمکی برای یافتن قدرتمند قابل مقایسه ای را فراهم میکند که اگر در برخی موارد برای فهرست کردن عالی نباشد، بنابراین می بینیم که بار دیگر اینجا چقدر خود کار سازی به طور چشمگیری هزینه ها را کاهش می دهد، چون می تواند بر روی یک مقیاس بزرگ برده شود (استک، 2005) (هال گریسون، 2006).

متأسفانه، خود کار سازی، محدودیت هایی دارد و بررسی دستی باید در برخی موارد انجام شود. برای مثال کشف می تواند به طور دستی با اتوماتیک انجام شود. وقتی کار به طور دستی انجام شود می تواند یک فعالیت خاص یا یک محصول فرعی از فعالیت های دیگر باشد مانند DACHS (لچر، 2006) و نمایش بایگانی های وب (وثرمان (2) و دیگران، 2002) Achipol این نوع رویکرد معمولاً برای بایگانی مرکزی عنوان در انجام می شود. اگر چه خزش عنوان به طور مؤثر برای کشف مطلب سایت های یا صفحه های مربوط آزمایش می گردد. ابزارهای خود کار می توانند به طور مطمئن (نه در این زمان) در مقایسه با یک شبکه متخصصان، مرجع هایی را برای بهترین مطالب که از آن مطلع هستند، فراهم نمایند.

به هر حال، فقدان حوزه تخصصی و عدم درک تنها معایب خزش گرها نیستند. این مسئله مورد توجه است که تأخیر برای یافتن سایت های جدید مورد نیاز است. زمان بسیار زیادی برای یافتن سایت های کلی گرفته می شود. وقتی وارد سایت های زودگذر (بی دوام) می شود، برای مثال مرتبط با یک رویداد، تاخیری می تواند بسیار طولانی باشد و آن ها را بایگانی کند این تفاوت به وسیله (میزانس، 2005) با یک مقایسه سایت های کشف شده توسط خزشگر الکسا و قابلیت دسترسی و روزانه به بایگانی اینترنت و سایت های مرتبط با انتخابات فرانسه در سال 2002 که توسط گروهی از کتابداران مرجع انجام شد و به وسیله کتابخانه ملی فرانسه بایگانی گردید، مورد بررسی قرار گرفته است.

این بررسی، مزیت های آشکاری را برای گزینش فعالیت دستی در مجموعه های مبتنی بر ربط، برای کشف به موقع و در تمرکز عمیق نشان می دهد.

رده بندی بایگانی وب بر طبق روش بررسی اش همچنین می توانست در یک شاخه مالی انجام شود. بر تراز پردازش اتوماتیک و دستی دو حالت، شخص می تواند برای مثال نوع منبع استفاده شده برای کشف، تناوب جست و جو، و ضبط، سطح کیفیت بازبینی ایجاد شده و دانه دانه بودن آیتم های بایگانی شده (سایت ها، صفحه ها)، و مانند آن را مورد بررسی قرار می دهد.

ص: 50

1- فیلیپز (2005) تخمین های بسیار مفید و دارای جزئیات در خصوص زمان و هزینه های فرآیند دستی خود کار سازی سایت ها برای یکی از قدیمی ترین آرشیوهای وب موجود را فراهم کرده است. تخمین های زمانی به صورت ذیل هستند: - تشخیص و انتخاب 30 دقیقه - گردآوری، اطمینان از کیفیت، و موارد آرشیوی: 210 دقیقه - فهرست نویسی 81 دقیقه

به هر حال، این امر که بیشتر بایگانی های وب به دو مدل عمده تمایل دارند، تمایز اصلی است که آیا انتخاب به طور دستی یا غیر آن انجام شده است. یکی مدل های خزش گر های کلی است که معمولاً مبتنی بر مرکز (دامنه های ملی یا دامنه های عمومی) است یا به صورت آزاد (بایگانی اینترنت)، و دیگری مدل گزینش فردی تعداد محدودی از جست و جوها یا نقاط ورودی است که به طور دستی انجام می شود (معمولاً سایت ها). تمایز در بیشتر در رویکرد روش شناختی آن ها به ندرت دیده شده یا هیچ یک برای طبقه بندی آن ها مورد استفاده قرار نمی گیرد.

6- نتیجه گیری

وب فقط 15 سال است که وجود دارد و می توان گفت که نگهداری و حفاظت از حافظه اش، به طور نسبی در مقایسه با رسانه های دیگر از ابتدا آغاز شده بود. (1) اما فقط مراحل لازم اولیه برای حفاظت آن ایجاد شده است. وضعیت حفاظت جاری، به تعداد بسیار کمی از مؤسسه ها وابسته است و پوشش زیادی را حاصل نشده است.

نقش ها و مسئولیت ها برای بیشتر سهامداران بسیار واضح و آشکار نیست و توانایی در حمایت از بسیاری از مجموعه های شاخص به وجود نیامده است. و ما هنوز در دوره ای هستیم که هیچ گسیختگی فناورانه از زمان آغاز وب، رخ نداده است.

مرورگرهای جاری با یکدیگر با تعداد محدودی که برنامه های کامپیوتری متصل می توانند فرمت های زیادی را جابه جا کنند که می تواند بر روی وب یافت شود.

اما این موقعیت، تا ابد طول نخواهد کشید و حفاظت از وب با یک چالش جدی روبه رو خواهد شد وقتی تغییرات مهم فناوری در وب رخ می دهد (که ممکن نیست مانند اینها بعداً دیده شود).

بنابراین، دلگرم کننده است که بینیم که بسیاری از مؤسسه های (حفظ) میراث، در بایگانی وب در حال به کارگیری هستند. بررسی اخیر توسط گروه پژوهش کتابخانه (RLG2006) نشان داد که 60 درصد اعضای مورد بررسی شان، بایگانی وب را قسمتی از مأموریت خود پنداشته اند (RLG2006) که بسیار دلگرم کننده است.

امیدواریم که بازنمون های ایجاد شده در این فصل، مباحث مهم و روش ها، متفقاً با منطق و دلیل، به آن ها و دیگران برای مشارکت در این تلاش گروهی، کمک خواهد کرد.

ص: 51

- Aarseth, E. J. (1997). *Cybertext: perspectives on ergodic literature*. Baltimore, MD: Johns Hopkins University Press . 1
- Abiteboul, S., Cobena, G., Masanès, J., Sedrati, G. (2002). A first experience in archiving the French Web. Paper presented at the Proceedings of the 6th European Conference on Research and Advanced Technology for Digital Libraries . 2
- Abiteboul, S., Preda, M., Cobena, G. (2003). Adaptive on-line page importance computation. Paper presented at the Proceedings of the twelfth international conference on World Wide Web . 3
- Antoniol, G., Canfora, G., Cimitile, A., De Lucia, A. (1999). Websites: files, programs or database. Paper presented at the 1st International Workshop on Web Site Evolution, Atlanta, USA . 4
- Arvidson, A., Persson, K., Mannerheim, J. (2000). The Kulturarw3 project – The Royal Swedish Web Archive – An example of "complete" collection of web pages. Paper presented at the 66th IFLA – International Federation of Library Associations and Institutions, Jerusalem . 5
- Baeza-Yates, R. Castillo, C. (2005). Characteristics of the Web of Spain. *Cybermetrics*, 9 . 6
- Baeza-Yates, R., Castillo, C., Efthimiadis, E. (2005a). Characterization of national Web domains . 7
- Baeza-Yates, R. A., Castillo, C., Marin, M., Rodriguez, A. (2005b). Crawling a country: better strategies than breadth-first for Web page ordering. Paper presented at the WWW 05: Proceedings of the 14th international conference on World Wide Web, Chiba, Japan . 8
- Balayé, S. (1988). *La Bibliothèque nationale, des origines a 1800 (Histoire des idées et critique littéraire; vol. 262)*. Genève: Droz . 9
- Battelle, J. (2005). Google Announces New Index Size, Shifts Focus from Counting. <http://battellemedia.com/archives/001889.php> . 10
- Benjamin, W. (1963). *Das Kunstwerk im Zeitalter seiner technischen Reproduzierbarkeit; drei Studien zur Kunstsoziologie*. [Frankfurt am Main]: Suhrkamp . 11
- Bergman, M. I. K. (2001). The deep Web: Surfacing hidden value. *The Journal of Electronic* . 13

- Bergmark, D. (2002). Collection synthesis. Paper presented at the 2nd ACM/IEEE-CS joint conference .14
on Digital libraries, Portland, USA
- Bergmark, D., Lagoze, C., Sbityakov, A. (2002). Focused crawls, tunneling, and digital libraries. Paper .15
presented at the 6th European Conference on Research and Advanced Technology for Digital Libraries,
Roma, Italy
- Berners-Lee, T. Connolly, D. (1995). Hypertext Markup Language – 2.0. RFC,1866 .16
- Berners-Lee, T. (1994). Universal Resource Identifiers in WWW, A Unifying Syntax for the Expression .17
of Names and Addresses of Objects on the Network as used in the World- Wide Web. RFC 1630
- Berners-Lee, T. (1998). Cool URIs don't change. [http://www.w3.org/Provider/ Style/ URI.html](http://www.w3.org/Provider/Style/URI.html) .18
- and ultimate destiny of the World Wide Web by its inventor (1st pbk. ed.). New York: HarperCollins .19
- .Björneborn, L. Ingwersen, P. (2001). Perspective of webometrics .20
Scientometrics, 50(1), 65–82 .21
- Bolter, J. D. (2001). Writing space: Computers, hypertext, and the remediation of print (2nd ed.). .22
Mahwah, NJ: Lawrence Erlbaum Associates
- Borgman, C. L. (2000). From Gutenberg to the global information infrastructure: access to information .23
in the networked world (Digital libraries and electronic publishing). Cambridge, MA: MIT
- Borgman, C. L. (2003). The Invisible Library: Paradox of the Global Information Infrastructure. Library .24
Trends, 51(4), 652–674
- Boudrez, P. Eynde, V. D., Sofie. (2002). Archiving Websites .25
- Boufkhad, Y. Viennot, L. (2003). The Observable Web. RR .26
- Boyko, A. (2004). Test Bed Taxonomy. IIPC Reports, 16 .27
- Broder, A., Kumar, R., Maghoul, F., Raghavan, P., Rajagopalan, S., Stata, R., et al. (2000). Graph .28
structure in the web. Paper presented at the 9th International World Wide Web Conference (WWW9),
Amsterdam, Netherlands

Brown, A. (2006). Archiving the Web: A guide for information management professionals. Library . 29
.Assn Pub

Brügger, N. (2005). Archiving Websites, general considerations and strategies. A arhus, Denmark: . 30
Center for Internet Research

ص: 53

- Bruns, A. (2005). Gatewatching: Collaborative online news production (Digital formations, v. 26).. 31
New York: P. Lang
- Burner, M. (1997). Crawling towards Eternity Building An Archive of The World Wide Web. New .32
Architect, 5
- Canfora, L. (1989). The vanished library (Hellenistic Culture and Society; 7). Berkeley: University of .33
California Press
- Canfora, L. (1996). Les bibliothèques anciennes et l'histoire des textes. In M. Baratin, C. Jacob (Eds.), .34
Le pouvoir des bibliothèques: la mémoire des livres en Occident. (pp. 338 p). Paris: A. Michel
- Carlin, J. W. (2004). Harvest of agency public websites. NARA Bulletin, 2005-02 .35
- Castells, M. (1996). The rise of the network society. Malden, MA: Blackwell .36
- Castillo, C., Marin, M., Rodriguez, A., Baeza-Yates, R. A. (2004). Scheduling Algorithms for Web .37
Crawling
- Chakrabarti, S. (2002). Mining the Web: discovering knowledge from hypertext data. San Francisco, .38
CA: Morgan Kaufmann
- Chakrabarti, S., Berg, M. V. D., Dom, B. (1999). Focused crawling: A new approach to topic-specific .39
Web resource discovery. Computer Networks (Amsterdam, Netherlands: 1999), 31, 1623-1640
- Chang, K. C.-C., He, B., Li, C., Patel, M., Zhang, Z. (2004). Structured .40
,databases on the web: observations and implications. SIGMOD Record .41
61-70 ,(3)33 .42
- Charlesworth, A. (2003). Legal issues relating to the archiving of Internet resources in the UK, EU, USA .43
and Australia
- Cho, J., Garcia-Molina, H. (2000). The evolution of the web and implications for an Incremental .44
Crawler. Paper presented at the Proceedings of the 26th International Conference on Very Large Data Bases
- Cho, J., Garcia-Molina, H., Page, L. (1998). Efficient Crawling Through url ordering. Computer .45
Networks and Isdn Systems, 30, 161-172

Christensen-Dalsgaard, B. (2001). Archive experience, not data. Paper presented at the Preserving the .46
Present for the Future – Strategies for the Internet, The Royal Library, Copenhagen, Denmark

Crowston, K., Williams, M. (1997). Reproduced and emergent genres of communication .47

ص: 54

- on the World-Wide Web. Paper presented at the 30th Annual Hawaii International Conference on System Sciences (HICSS-30), Wailea, USA
- Cruse, P., Eckman, C., Kunze, J. (2003). Web-based government information: Evaluating solutions for .48 capture, curation, and preservation. An Andrew W. Mellon funded initiative of the California Digital Library
- Dahn, M. (2000). Counting Angels on a Pinhead: Critically Interpreting Web Size Estimates. Online, .49
January/February, 35-40
- Day, M. (2006). The long-term preservation of Web content. In J. Masanè s (Ed.), Web archiving. .50
Berlin Heidelberg New York: Springer
- Dikaiakos, M. D. (2004). Intermediary infrastructures for the World Wide web. Computer Networks, .51
45(4), 421-47
- Dobra, A., Fienberg, S. E. (2004). How Large Is the WorldWide Web?. In M. Levene, A. Poulouvassilis .52
(Eds.), Web dynamics web dynamics – adapting to change in content, size, topology and use. (pp. 23-44).
Berlin Heidelberg New York: Springer
- Dubberly, H., Forlizzi, J., Hodge, C., Laurel, B., Lyman, P., Meggs, P. B., et al. (2002). Archiving .53
experience design, a virtual roundtable discussion. LOOP: AIGA Journal of Interaction Design Education,
Number 6
- Dumais, S. T., Cutrell, E., Cadiz, J. J., Jancke, G., Sarin, R., et al. (2003). Stuff I've seen: A system for .54
personal information retrieval and re-use. Toronto, Canada
- Egghe, L. (2000). New informetric aspects of the Internet: some reflections – many problems. Journal of .55
Information Science, 26(5), 329-335
- Eisenstein, E. L. (1979). The printing press as an agent of change: Communications and cultural .56
transformations in early modern Europe. Cambridge [Eng.]; New York: Cambridge University Press
- (Entlich, R. (2004). Blog Today, Gone Tomorrow? Preservation of Weblogs. RLG DigiNews, 8(4 .57
- Eriksen, L. B. Ihlström, C. (2000). Evolution of the web news genre – The slow move beyond the print .58
metaphor. Paper presented at the 33rd Hawaii International Conference on System Sciences (HICSS-33),
Hawaii, USA

Estivals, R. (1961). Le dé pôt lé gal sous l'Ancien Ré gime, de 1537 a 1791. Paris: M. Rivière .59

Estivals, R. (1965). La statistique bibliographique de la France sous la monarchie au .60

ص: 55

- Fauconnier, S. Frommé, R. (2004). Capturing unstable media, summary of research .61
- Fayet-Scribe, S. (2000). Histoire de la documentation en France: Culture, science, et technologie de l'information, 1895-1937 (CNRS histoire). Paris: CNRS .62
- Featherstone, M. (2000). Archiving cultures. *British Journal of Sociology*, 51(1) .63
- Febvre, L.P.V. Martin, H. J. (1976). The coming of the book: The impact of printing 1450-1800 ([New ed.] ed.). London: NLB .64
- Fielding, R. T., Gettys, J., Mogul, J., Nielsen, H. F., Masinter, L., J, P., et al. (1999). Hypertext Transfer Protocol - HTTP/1.1. RFC, 2616 .65
- Fitch, K. (2003). Web site archiving: An approach to recording every materially different response produced by a website. Paper presented at the AusWeb (2003): The Ninth Australina World Wide Web Conference, Sanctuary Cove, Australia .66
- Florescu, D., Levy, A., Mendelzon, A. (1998). Database techniques for the World- Wide Web: A survey. *SIGMOD Record* 27, 59-74 .67
- Freeman, E. Gelernter, D. (1996). Lifestreams: A storage model for personal data. *SIGMOD Record*, 25(1), 80-86 .68
- Gemmell, J., Bell, G., Lueder, R., Drucker, S., Wong, C. (2002). MyLifeBits: fulfilling the Memex vision. Juan-les-Pins, France .69
- Gibson, D., Punera, K., Tomkins, A. (2005). The volume and evolution of web page templates. Paper presented at the WWW'05 14th international conference on World Wide Web, Chiba, Japan .70
- Gillies, J. Cailliau, R. (2000). How the Web was born: The story of the World Wide Web. Oxford: Oxford University Press .71
- Golder, S. Huberman, B. A. (2005). The Structure of Collaborative Tagging Systems 73. Gomes, D. .72
- Silva, M. J. (2003). A Characterization of the Portuguese Web. Paper presented at the 3rd Workshop on Web Archives (IWA'03), Trondheim, Norway .73

Gulli, A. Signorini, A. (2005). The indexable web is more than 11.5 billion pages. Chiba, Japan .74

Halavais, A. (2004). Tracking Ideas in the Blogosphere .75

Hallgrímsson, T. (2006). Access and finding aids or web archives. In J. Masanè s (Ed.), Web archiving. .76
Berlin Heidelberg New York: Springer

Hine, C. (2000). Virtual ethnography. London; Thousand Oaks, CA: Sage .77

ص: 56

- Hofmann, M. Beaumont, L. R. (2005). Content networking: Architecture, protocols, and practice (The Morgan Kaufmann Series in Networking). Amsterdam; Boston: Morgan Kaufmann
- Ingwersen, P. (1998). The calculation of web impact factors. *Journal of Documentation*, 54(2) .79
- Jones, S. Johnson, C. (2006). Web Use and Web Studies. In J. Masanè's (Ed.), *Web Archiving*. Berlin Heidelberg New York: Springer .81
- Jones, W., Bruce, H., Dumais, S. (2001). Keeping found things found on the web. Atlanta, GA, USA .82
- Jones, W., Bruce, H., Dumais, S. (2003). How do people get back to information on the Web? How can they do it better? Paper presented at the IFIP INTERACT'03 .83
- Kahle, B. (1997). Preserving the Internet. *Scientific American*, 397, 82-84 .85
- Kahle, B. (2002). The Internet Archive. *RLG DigiNews*, 6(3) .86
- Kimpton, M., Braggs, M., Ubois, J. (2006). Year by Year: From an Archive of the Internet to an Archive on the Internet. In J. Masanè's (Ed.), *Web Archiving*. Berlin Heidelberg New York: Springer .87
- Koehler, W. (1999). Unraveling the ISSUES, ACTORS, ALPHABET SOUP of the Great Domain Name Debates. *Searcher*, 7(5) .88
- Koehler, W. (2004). A longitudinal study of Web pages continued: a consideration of document persistence. *Information Research*, 9(2) .89
- Krishnamurthy, B. Rexford, J. (2001). Web protocols and practice: HTTP/1.1, networking protocols, caching, and traffic measurement. Boston, MA: Addison-Wesley 91.
- Lagoze, C., Dean B. K., Sandy, P., Jesurogaili, S. (2005). What Is a Digital Library Anymore, Anyway? *Beyond Search and Access in the NSDL. D-Lib Magazine*, 11-11 92.
- Lamos, C., Eirinaki, M., Jevtuchova, D., Vazirgiannis, M. (2004). Archiving the Greek Web. Paper presented at the 4th International Web Archiving Workshop (IWA'04), (Bath (UK) .93
- Landow, G. P. (1997). *Hypertext 2.0* (Rev., amplified ed.). Baltimore: Johns Hopkins University Press .93
- Lavoie, B. F. Schonfeld, R. C. (2005). The systemwide print book collection. Paper presented at the CNI .94

- Lawrence, S. Giles, C. L. (1998). Searching the Web. *Science*, 281, 175 .95
- Lawrence, S. Giles, C. L. (1999). Accessibility of Information on the Web. *Nature*, 400, 107–109 .96
- Lecher, H. E. (2004). Informant networks, alarm systems, and research contributors. Selection and ingest process for the Digital Archive for Chinese Studies. Paper presented at the Archiving Web Resources Conference – Issues for Cultural Heritage Institutions, NLA, Canberra, Australia .97
- Lecher, H. E. (2006). Academic Web archiving: DACHS. In J. Masanè s (Ed.), *Web archiving*. Berlin Heidelberg New York: Springer .98
- Levy, P. (1997). *Collective intelligence: Mankind's emerging world in cyberspace*. Cambridge, MA: Perseus Books .99
- Liu, C. Albitz, P. (1999). *DNS BIND (3rd ed.)*. O'Reilly Associates .100
- Lueg, C. Fisher, D. (2003). *From Usenet to CoWebs: Interacting with social information spaces (Computer supported cooperative work)*. Berlin Heidelberg London New York: Springer .101
- Lyle, J. A. (2004). Sampling the Umich.edu Domain. Paper presented at the 4th International Web Archiving Workshop (IWAW'04), Bath (UK) .102
- Lyman, P. (2002). Archiving the World Wide Web. In CLIR (Ed.), *Building a national strategy for preservation: issues in digital media archiving*. Council on Library and Information Resources and the Library of Congress .103
- Lyman, P. Kahle, B. (1998). Archiving digital cultural artifacts. *D-Lib Magazine* .105
- Mantratzis, C. Orgun, M. (2004). Towards a peer2peer world-wide-web for the broadband-enabled user community .106
- Masanè s, J. (2002). Towards continuous Web archiving: First results and an agenda for the future. *D-Lib Magazine*, 8(12) .107
- Masanè s, J. (2004). Site-first priority: Implementing the frontline .108
- Masanè s, J. (2005). Web archiving methods and approaches: A comparative study. *Library Trends* .109

Masanè's, J. (2006a). Collecting the hidden web. In J. Masanè's (Ed.), Web archiving. Berlin . 110
Heidelberg New York: Springer

Masanè's, J. (2006b). Selection for Web Archives. In J. Masanè's (Ed.), Web archiving. Berlin . 111
Heidelberg New York: Springer

- Mohr, G., Kimpton, M., Stack, M., Ranitovic, I. (2004). Introduction to Heritrix, an archival quality web crawler. Paper presented at the 4th International Web Archiving Workshop (IWAW'04), Bath (UK) .112
- Mueller, M. (2002). Ruling the root: Internet governance and the taming of cyberspace. Cambridge, MA: MIT .114
- Najork, M. Heydon, A. (2001). High-performance Web crawling. SRC Research Report .115
- Najork, M. Wiener, J. (2001). Breadth-first search crawling yields high-quality pages. Paper presented at the 10th World Wide Web Conference (WWW'10), Hong Kong .116
- National Archives of Australia. (2001). Archiving Web resources: A policy for keeping records of web-based activity in the Commonwealth Government .117
- Osborn, T. (1999). The ordinariness of the archive. *History of the human sciences*, 12(2) .118
- Page, L., Brin, S., Motwani, R. Winograd, T. (1998). The Pagerank citation ranking: Bringing order to the Web, 17 .119
- Pandey, S. Olston, C. (2005). User-centric Web crawling. Chiba, Japan .120
- Pant, G., Srinivasan, P. Menczer, F. (2004). Crawling the Web. In M. Levene, A. Poulovassilis (Eds.), *Web Dynamics*. (pp. 153-178). Berlin Heidelberg New York: Springer .121
- Pastor-Satorras, R. Vespignani, A. (2004). Evolution and structure of the Internet: A statistical physics approach. Cambridge, UK; New York: Cambridge University Press .122
- Phillips, M. E. (2005). Selective archiving of Web Resources: A study of acquisition costs at the National Library of Australia. *RLG DigiNews*, 9(3) .123
- Qin, J., Zhou, Y. Chau, M. (2004). Building domain-specific web collections for scientific digital libraries: A meta-search enhanced focused crawling method. Tuscon, AZ, USA .124
- Rekimoto, J. (1999). Time-machine computing: A time-centric approach for the information environment. Paper presented at the 12th annual ACM symposium on User interface software and technology, Asheville, North Carolina, USA .126

Riché , P. (1996). La bibliothè que et la formation de la culture mé dié vale. In M. Baratin, C. Jacob) .127
(Eds.), Le pouvoir des bibliothè ques: la mé moire des livres en Occident (p

ص: 59

- Ringel, M., Cutrell, E., Dumais, S., Horvitz, E. (2003). Milestones in Time: The Value of Landmarks . 128
in Retrieving Information from Personal Stores. Paper presented at the IFIP INTERACT '03
- ?RLG. (2006). Web Archiving Program. <http://www.rlg.org/en/page.php> .129
- Page-ID=399 .130
- Roche, X. (2006). Copying web sites. In J. Masanè s (Ed.), Web Archiving. Berlin Heidelberg New . 131
York: Springer
- Rosenfeld, L. Morville, P. (2002). Information architecture for the World Wide Web (2nd ed.). . 132
Cambridge, MA: O'Reilly
- Scharl, A. (2000). Evolutionary Web development (Applied computing). Berlin Heidelberg New . 133
York: Springer
- Shepherd, M. Polanyi, L. (2000). Genre in Digital Documents. Paper presented at the Proceedings of . 134
the 33rd Hawaii International Conference on System Sciences - vol. 3
- Sonnenreich, W. (1997). A History of Search Engines. [http://www.wiley.com/legacy/](http://www.wiley.com/legacy/compbooks/sonnenreich/history.html) . 135
[compbooks/sonnenreich/history.html](http://www.wiley.com/legacy/compbooks/sonnenreich/history.html)
- Spinellis, D. (2003). The decay and failures of web references. Communications of ACM, 46(1), 71- . 136
77
- Stack, M. (2005). Full Text Search of Web Archive Collections. Paper presented at the IWWAW'05, . 137
Vienna, Austria
- Star, S. L. Ruhleder, K. (1994). Steps towards an ecology of infrastructure: Complex problems in . 138
design and access for large-scale collaborative systems. Chapel Hill, NC, United States
- Teevan, J. (2004). How people re-find Information when the Web changes. AIM- 2004-012 .139
- Thelwall, M. (2001). Extracting macroscopic information from Web links. Journal of the American . 140
Society for Information Science and Technology, 52(13), 1157-1168 141. Thelwall, M. (2006).
Interpreting social science link analysis research: A theoretical framework. Journal of American Society of

Thelwall, M. Harries, G. (2004). Do the websites of higher rated scholars have .142

ص: 60

significantly more online impact? Journal of the American Society for Information Science and Technology,
55(2), 149-59

Thelwall, M. Vaughan, L. (2004). A fair history of the Web? Examining country balance in the Internet .143
archive. Library Information Science Research, 26(2), 162- 176

Ubois, J. (2002). The Oakland archive policy. Recommendations for managing removal requests and .144
preserving archival integrity

Voerman, G., Keyzer, A., Hollander, F. D., Druiven, H. (2002). Archiving the Web: Political Party) .145
Web sites in the Netherlands. European Political Science, 2 (1

پیشرفت در رایانه و ارتباطات این امکان را فراهم ساخته که ذخیره کتاب، صفحه موسیقی، فیلم، بسته های نرم افزاری و تمام صفحات عمومی وب که تاکنون ساخته شده اند هزینه - سودمندی داشته باشند و دسترسی به این مجموعه ها از طریق اینترنت برای همه، از جوان تا پیر، در سرتاسر دنیا فراهم شود. برای سال های آینده رسالت آرشیو مشخص است: ایجاد کتابخانه ای جهانی که تمام دانش به سهولت در اختیار هر مرد و زن و کودک در سرتاسر جهان قرار گیرد. در مقاله حاضر بحث آرشیو اینترنت و آرشیو در محیط اینترنت، دسترس پذیری دراز مدت تمام دانش برای همه افراد جهان پیشرفت های فن آورانه ایجاد آرشیو اینترنتی مورد بررسی قرار می گیرد. سپس برنامه ها و پروژه های نمونه آرشیو اینترنتی در دنیا آرشیو اروپا و پتاباکس معرفی می شوند.

*از آرشیو اینترنت تا آرشیو در اینترنت (1)

میشل کیمتون (2) | جف یوبویس (3)

ترجمه: مرضیه هدایت (4)

مقدمه

«آرشیو اینترنت» (5)، از آغاز کار خود در 1995، تاکنون، هدف درازمدت فراهم کردن دسترسی جهانی به تمام دانش در طول عمر را دنبال کرده است.

در 10 سال گذشته، آرشیو، در برنامه های کشورهای مختلف، برای کاربران متعدد و گوناگون، شرکت داشته و برای حفظ صفحات وبی، کتاب، موسیقی، نرم افزار، و تصاویر پویا فعالیت کرده تا آن ها را از طریق اینترنت دسترس پذیر سازد.

در حال حاضر، ماشین Wayback آرشیو، روزانه به 70,000 بازدید کننده، و در هر ثانیه به 200 درخواست پاسخ می دهد. مجموعه 600 ترابایتی آن شامل 50 بیلیون صفحه وب، 30000 جلد کتاب 36000 قطعه موسیقی، و 15000 قطعه فیلم است.

حضور صدها مشترک این امر را امکان پذیر ساخته است. به واسطه پیشرفت های فنی در صنعت رایانه،

ص: 63

Year - by - Year: From an Archive of the Internet to an Archive on the Internet: in Masanes, Julien (ed.), - 1
.Web Archiving. Berlin Heidelberg New York: Springer.pp.201-212

Michele Kimpton -2

Jeff Ubois -3

4- عضو هیئت علمی سازمان اسناد و کتابخانه ملی ایران

Internet Archive -5

امکانات این سیستم هر 12 - 18 ماه دو برابر می شود. در جامعه آرشیوی 10 سال زمان زیادی نیست، اما زمان مناسب برای خلق یک عامل 100 برابر پیشرفت است: از سال 1997 تاکنون قیمت دیسک خام بیش از 99 درصد کاهش یافته است؛ یعنی از 180 دلار به 50 سنت برای هر گیگابایت رسیده و بیش از 25 میلیون ارتباط پهنای باند فقط در ایالات متحده اضافه شده است.

با توجه به اینکه میزان پیشرفت رایانه بر اساس ذخیره و بازیابی است، می توان آن چه را که راج ردی (1) از دانشگاه کارنگی ملون (2) با عبارت «دستیابی جهانی به تمام دانش» بیان کرده است، محتمل دانست.

پیشینه: چاپ اینترنتی اولیه

از اواسط دهه 1980 معلوم بود که تغییری در دنیای چاپ الکترونیکی در حال ظهور است. در نیمه اول دهه 1990، حرکت روزنامه ها از شکل پایگاه های بسته با مالکیت انحصاری به شکل اینترنتی شروع شد و نظریه اینترنت به عنوان کتابخانه شروع به شکل گیری کرد. نظام های چاپی اینترنتی از قبیل (3) WAIS و Gopher، به عنوان متمم های صفحات وبی دیده شدند؛ استعاره اینترنت به عنوان کتابی با صفحات وبی که فهرست مندرجات آن توسط سرورهای Gopher و نمایه آن توسط سرورهای WAIS تهیه می شد، به عنوان جایگزینی برای استعاره اینترنت به عنوان «فرشاهراه» محسوب شد.

شروع به کار خدمات آلتا ویستا (4) در دسامبر 1995 ثابت کرد که تمام صفحات موجود در وب را می توان مجموعه ای یگانه دانست که قابل نمایه شدن و جست و جو کردن بر روی شبکه برای تمام کاربران است؛ اما معلوم نبود که این صفحات را چگونه باید حفظ و نگهداری کرد.

آغاز به کار آرشیو اینترنت

آرشیو اینترنت، به طور رسمی با همکاری بروس گیلیت (5) و بروستر کال (6)، در آوریل 1996 شروع شد. در آن زمان، پیوندهای شکسته (404 درگاه) یک مشکل روبه رشد بود و روشن بود که بیشتر صفحات وبی عمر کوتاهی دارند. برای این مشکل راه حلی مورد نیاز بود و نظامی برای آرشیو صفحات وبی قبل از اینکه پاک شوند یک راهکار ضروری به نظر می رسید.

این مسئله، به تصمیم گیری راجع به طرح اولیه در «آرشیو» درباره سیاستگذاری مجموعه ها منجر شد؛ مانند حریص بودن نسبت به جمع آوری مطالبی که در معرض خطر از بین رفتن بودند، و شکار فرصت ها برای جمع آوری و رقومی کردن مواردی از گذشته مثل پستینگ های یوزنت (7).

با این حال، هنوز در 1996، در جامعه اینترنتی اهمیت از دست رفتن صفحات وبی مشکل چندان حساسیت

ص: 64

Raj Reddy -1

Carnegie Mellon University -2

(Wide Area Information Servers (WAIS -3

Alta Vista -4

Bruce Gilliat -5

Brewster Kahle -6

Usenet -7

برانگیز نبود. چون وب سابقه تاریخی چندانی نداشت توضیح فواید صفحات از دست رفته مشکل بود.

برای نشان دادن ارزش بالقوه چنین صفحاتی، آرشیو، با مؤسسه اسمیتسونین (1) در واشنگتن دی.سی، برای جمع آوری نسخه فوری (2) وبگاه های تمام نامزدهای ریاست جمهوری 1996 (3) همکاری کرد. ابزار انجام این پروژه چندان مورد رضایت نبود. این ابزار به طور اساسی عبارت بودند از: تجهیزات رایانه های شخصی که بر پایه ضبط وب سایت ها از طریق دنبال کردن پیوندها از صفحه اصلی عمل می کردند

این داده ها، در نهایت، به آرشیو ریاست جمهوری اسمیتسونین منتقل شد که در حال حاضر شامل

صفحاتی از 5 حزب سیاسی و کاندیداهای پرشمار ریاست جمهوری از بیل کلینتون گرفته تا پت بوکانن (4) است. بسیاری از سایت ها در این آرشیو با حذف کاندیدها بسته شدند.

بر اساس این موفقیت، کتابخانه کنگره به آرشیو مأموریت داد تا مجموعه پیوسته متمرکزی از انتخابات سال 2000 ایجاد کند و این درخواست را نیز برای انتخابات 2002 تجدید کرد.

همچنین، در سال 1996، آرشیو، ارتباطش را با اینترنت الکسا شروع کرد. اینترنت الکسا، مؤسسه ای انتفاعی است که خزشگری و آرشیوسازی وب را در نوامبر شروع کرد تا داده های پیشنهادی نوار ابزار مرورگر درباره سایت های دیده شده را Plug-in (وصلینه) کند و بر اساس داده های جمع آوری شده از سایر کاربران و پیشنهاد درباره صفحات مرتبط - که ممکن است مورد توجه باشند - عمل می کرد.

دو پیشرفت دیگر از سال 1996 نیز قابل ذکرند. نخستین پیشرفت، فناوریانه است. در 1996، هنوز ارجحیت نوار نسبت به دیسک از نظر قیمت قابل توجه بود و آرشیو اولین نسل زیرساختار را با استفاده از روپات های ذخیره نواری ساخت که با ADIC 50 شروع شد. با وجود همکاری های سخاوتمندانه فروشندگان اصلی، در نهایت ثابت شد این کار قابل دفاع نیست. نیازهای دستیابی که توسط کاربران آرشیو مطرح می شدند بیش از حد شدید و زمان بازایی خیلی کند بود. همانگونه که بروس گیلیت به طنز می گوید: گرفتن یک صفحه می تواند در چند ثانیه باشد... یا چند روز بعد.

دومین پیشرفت به مسائل قانونی مربوط می شد. مسائل قانونی جمع آوری صفحات وبی با الگوریتم حریمانه و ارائه آن ها بر پایه «صرف نظر کردن» (5) آن گونه که آرشیو در آن سال شروع به انجامش کرد - روشن نبود. بهبود پروتکل ارتقای روپات های محدودیت متنی متعلق به آلتا ویستا گامی مهم بود چون «پیش فرض ها» را تغییر داد - صاحبان صفحات وبی که از گذاشتن آن ها در نمایه موتور جست و جو یا در آرشیو ابا داشتند، روشی ساده برای صرف نظر کردن یا انتقال صفحاتی که صاحب آن بودند (اثبات آن از طریق توان آن ها برای مناسب سازی راهنمای پایه ای وبگاه مورد نظر بود) داشتند. روپات های متنی برای انتقال موارد مورد نظر برای آن ها که صاحب صفحات در آرشیو بودند راه حلی ارائه داد، ولی نتوانست مشکل انتقال صفحات متعلق به دیگران را رفع کند. این موارد انتقالی، در کنفرانس کوچکی

ص: 65

Snapshots -2

<http://movie0.archive.org/96-Elections/index.htm> -3

Pat Buchanan -4

Opt-out -5

توضیح داده شد، که توسط آرشیو اینترنت در 2002، در یوسی برکلی (1) برگزار شد (یوبویس 2002).

ساختار پیوندی و روبات های نواری

جمع آوری صفحات وبی، داده های پیوندی و «رد پایهای استفاده ای (2)» به عنوان نمونه، فرصت هایی که میلیون ها کاربر وب با رفتن از صفحه ای به صفحه ای ایجاد کرده بودند؛ در اوایل 1997 با معرفی نوار ابزار الکسا، یک مرورگر Plug-in از طریق تهیه اطلاعات روی سایتی که در حال دیده شدن بود و نیز پیشنهادهایی برای سایت های مرتبط به کاربران کمک می کرد تا به ناوبری در وب پردازند.

داده های پیوندی و ردپاهای استفاده ای که توسط الکسا جمع آوری می شدند، به عنوان نظام فیلترینگ اشتراکی عمل می کرد و صفحاتی را برجسته می کرد که جامعه اینترنتی دارای بیشترین ارزش می دانست. پیوندها و کلیک ها بر اساس ارزش صفحه داده شده امتیاز می گرفتند.

توان مشخص کردن ارزش یک صفحه، به طور خودکار، تا حد زیادی با یک طرح اساسی دیگر مرتبط بود. در 1997، تعدادی از بزرگ ترین و موفق ترین کتابخانه های وبی، فهرست سایت هایی مانند یاهو بودند. اما در مورد نحوه سنجش راهکارهای فهرست نویسی دستنامه ای که ممکن بود زمانش بگذرد، اطمینانی نبود. آیا احتمال حذف فهرست نویسی دستنامه ای و پذیرش روند انتخاب های قطعی با دیدگاه «جمع آوری کامل» که با فراداده ایجاد شده توسط کاربر به جای فهرست ترکیب می شد، امکان پذیر بود؟

به نظر می رسید پاسخ مثبت باشد. به همین دلیل، الکسا بر اساس ثبت میزان استفاده، به خزش در صفحات پرداخت. سنجش آن از طریق داده های جمع آوری شده توسط نوار ابزار مرورگر انجام می شد. صفحاتی که بسیار دیده شده بودند، نخستین هایی بودند که برایشان نسخه پشتیبان تهیه می شد.

خزشگر الکسا، طوری تنظیم شده بود که هر 8 هفته یک نسخه فوری از وب تهیه کند، و هنوز این برنامه اجرا می شود، اگر چه اندازه هر خزش از یک ترابایت در سال 1997 به 100 ترابایت در 2004 رسیده است.

مورد دیگری که آرشیو در سال 1997 با آن روبه رو شد این بود که آیا به ذخیره روی نوار می توان اعتماد کرد یا ذخیره روی دیسک؟ از نظر قیمت هنوز نوار بر دیسک ارجح بود، ولی دستیابی کند بود.

همانگونه که در مقاله گری و شنوی (1999) (3) آمده «نسبت قیمت نوار، دیسک و RAM تقریباً 10:1000 است، به این معنی که ذخیره روی دیسک 10 بار گران تر از ذخیره روی نوار و ذخیره روی RAM 100 بار گران تر از ذخیره روی دیسک است».

اما وقتی هزینه دستیابی را با یارد آهنی بسنجیم، دیسک واقعاً بسیار ارزان تر است؛ «هزینه آرشیو نواری در مقایسه با هر ترابایت ذخیره روی دیسک نصف است، ولی دستیابی آسان به داده ها از طریق نوار ممکن نیست. هزینه هر دستیابی از طریق نوار، به صورت تصادفی، حدود 100 هزار بار بیشتر است (دستیابی از طریق دیسک: 100 دستیابی در ثانیه با هزینه 1 دلار در مقابل 10/000 دستیابی در ثانیه از

طریق نوار با هزینه 10000 دلار است)» (گری و شنوی 1999).

UC Berkeley -1

Usage trails -2

Gray and Shenoy -3

1998: حضور داده های آرشیوی بر روی (تقریباً) هر دسکتاپ

پس از آن که در 1996، نت اسکپ (1) در دسترس عموم قرار گرفت؛ در 1998، اینترنت و مشاغل وابسته به آن در یک تغییر جهانی در اولویت های سرمایه گذاری مورد توجه قرار گرفت. همانطور که بیلین ها دلار در بازارهای عمومی، سرمایه گذاری های مخاطره آمیز و راه اندازی اینترنت هزینه می شد، شماره صفحات در آرشیو هر 3 - 6 ماه دو برابر و تعداد کاربران اینترنت نیز هر چند ماه دو برابر می شد. به تدریج دستیابی به یک دغدغه تبدیل می شد.

در تلاش برای ایجاد دسترسی به موجودی و تأسیس آرشیو و اینترنت الکسا، به عنوان بخشی از زیر ساختار اینترنت، الکسا قراردادهایی با مرورگر های مایکروسافت و نت اسکپ بست. این امر به معنای حضور الکسا در 90 درصد رایانه های دنیا بود و کاربران چه می دانستند چه نه، الکسا بود که دسترسی به داده های تهیه شده توسط آرشیو اینترنت را به آن ها می داد.

برای آرشیو، ارائه داده ها برای ده ها میلیون کاربر تأثیر شدیدی بر زیر ساختار نوار مدار آن داشت.

در پایان 1998 دو امر مشخص بود:

- به دلیل میزان زیاد تقاضا رفتن از نوار به دیسک نیاز مبرمی بود. همان طور که تقاضای دسترس افزایش می یافت توان رویت های نواری برای پاسخگویی مورد تردید قرار می گرفت.

- سیاستگذاری های مجموعه های دستی گران تر از دیسک بود، یعنی آرشیو سازی بر پایه فهرست نویسی دستی، به ویژه وبگاه ها، بسیار گران تر از آرشیو سازی تمام سایت های قابل دسترسی بر اساس داده های جمع آوری شده از کاربران نهایی توسط الکسا بود.

1999: از نوار تا دیسک، یک خزشگر جدید و تصاویر متحرک

موفقیت تجاری خدمات الکسا، که به واسطه حضورش در هر رایانه شخصی مرتبط با اینترنت، بخش مشخصی بود، باعث شد آمازون، در سال 1999، الکسا را خریداری کند. این امر، در نهایت، به تغییراتی در ساختار هر دو سازمان منجر شد.

یک پیشرفت فناورانه مهم در سال 1999، ابداع خزشگر جدید توسط اندی جونل (2) بود. این خزشگر جدید، برای جمع آوری داده های وبی توانایی بیشتری داشت و از طریق ماشین های چندگانه قابل مدیریت بود. این خزشگر، به افزایش توانایی الکسا در فیلتر کردن 16 بیلین URL و 4 بیلین خزشگری و گسترش وسعت و عمق خزش خود پردازد.

همچنین، تصمیماتی درباره فرمت فایل ARC اخذ شد، که برای ذخیره صفحات وبی استفاده می شد. ویژه سازی فرمت فایل ARC - که برای نیازهای متعددی طراحی شده بود - از آغاز، در سال 1996 توسط مایک برنر (3) و بروستر کال توسعه یافت (برنر و کال 1996). نیازهای مذکور عبارت اند از:

● فایل باید خود شمول باشد: فایل باید اجازه بدهد اشیای متراکم، بدون استفاده از فایل نمایه همراه،

Netscape -1

Andy Jewel -2

Mike Burner -3

• فرمت باید برای پذیرش فایل های بازیابی شده از طریق پروتکل های شبکه ای متنوع مانند FTP ، HTTP ، اخبار، گوفر، و پست الکترونیکی قابل گسترش باشد؛

• فایل باید قابلیت «جریان داشتن» (1) داشته باشد: فایل باید بتواند فایل های چندگانه آرشیو را در جریان

داده ها به هم پیوند دهد؛

• رکورد باید با یکبار نوشتن ماندگار: باشد یکپارچگی فایل نباید به ایجاد بعدی نمایه درون فایلی محتوا وابسته باشد.

در عین حال که خاص سازی به نمایه بیرونی محتوا و شیء - جبرانی نیاز ندارد، چنین نمایه ای به میزان زیادی قابلیت بازیابی اشیای ذخیره شده در این فرمت را افزایش می دهد. در حال حاضر، آرشیو چنین نمایه هایی را فراهم کرده و به دنبال استاندارد سازی فرمت های آن ها از طریق کنسرسیوم حفاظت بین المللی اینترنتی (2) است.

در 999، حرکت از صفحات وبی، به سایر انواع داده ها نیز شروع شد. در این سال، هزینه ذخیره آن قدر کاهش یافت که آرشیو توانست به گردآوری تصاویر متحرک پردازد. از طریق شراکت با ریک پری لینگر (3)، از آرشیو Prelinger، پروژه دیجیتالی کردن 1000 فیلم (با هزینه حداکثر 160000 دلار) و نیز آرشیو کردن اخبار تلویزیون در پایان سال شروع شد.

2000: ایجاد مجموعه های موضوعی وب

در سال 2000، آرشیو به سطحی از ثبات فناوری رسید. پذیرش داده های خزشگری روش معمول شد و مهاجرت از نوار به دیسک سپری مجموعه گشت.

شناسایی و باز شوند؛

- فرمت باید برای پذیرش فایل‌های بازبایی شده از طریق پروتکل‌های شبکه‌ای متنوع مانند HTTP، FTP، اخبار، گوفر، و پست الکترونیکی قابل گسترش باشد؛
- فایل باید قابلیت «جریان داشتن»^۱ داشته باشد: فایل باید بتواند فایل‌های چندگانه آرشیو را در جریان داده‌ها به هم پیوند دهد؛
- رکورد باید با یکبار نوشتن ماندگار باشد: یکپارچگی فایل نباید به ایجاد بعدی نمایه درون فایلی محتوا وابسته باشد.

در عین حال که خاص‌سازی به نمایه بیرونی محتوا و شیء - جبرانی نیاز ندارد، چنین نمایه‌ای به میزان زیادی قابلیت بازبایی اشیای ذخیره شده در این فرمت را افزایش می‌دهد. در حال حاضر، آرشیو، چنین نمایه‌هایی را فراهم کرده و به دنبال استانداردسازی فرمت‌های آنها از طریق کنسرسیوم حفاظت بین‌المللی اینترنتی^۲ است.

در ۹۹۹، حرکت از صفحات وبی، به سایر انواع داده‌ها نیز شروع شد. در این سال، هزینه ذخیره آنقدر کاهش یافت که آرشیو توانست به گردآوری تصاویر متحرک بپردازد. از طریق شراکت با ریک پری لینگر^۳، از آرشیو Prelinger، پروژه دیجیتالی کردن ۱۰۰۰ فیلم (با هزینه حداکثر ۱۶۰۰۰۰ دلار) و نیز آرشیو کردن اخبار تلویزیون در پایان سال شروع شد.

۲۰۰۰: ایجاد مجموعه‌های موضوعی وب

در سال ۲۰۰۰، آرشیو به سطحی از ثبات فناوری رسید. پذیرش داده‌های خزشگری روش معمول شد و مهاجرت از نوار به دیسک سپری گشت.

جدول ۱. مجموعه‌های آرشیو اینترنت در مارس ۲۰۰۰

مجموعه	واحد	اندازه
وب (۱۹۹۶ تا ماه ۳ سال ۲۰۰۰)	۱ بیلیون صفحه	۱۳/۸ ترابایت
FTP (۱۹۹۶)	۵۰/۰۰۰ سایت	۰/۰۵ ترابایت
Usenet (۱۹۹۶ - ۱۹۹۸)	۱۶ میلیون پستینگ	۰/۵۹۲ ترابایت

در سال ۲۰۰۰، انتخابات دیگری در ایالات متحده برگزار شد و این بار تمام انتخاب‌کنندگان دسترسی اینترنتی داشتند. از نظر سیاسی، روشن بود که حضور در اینترنت برای برنده شدن حیاتی است و این امر تمرکز بر برخط بودن سیاسی را افزایش داد. آرشیو اینترنت، با کتابخانه کنگره برای گردآوری سایت‌های سیاسی شریک شد.

این اقدام، اولین پروژه آرشیو با کتابخانه کنگره بود و برای بسیاری از کارکنان آن به معنی انتقال از

1. Streamable

2. International Internet Preservation Consortium (IIPC) <http://netpreserve.org>

3. Rick Prelinger

جدول 1. مجموعه‌های آرشیوی اینترنت در مارس 2000

در سال 2000، انتخابات دیگری در ایالات متحده برگزار شد و این بار تمام انتخاب‌کنندگان دسترسی اینترنتی داشتند. از نظر سیاسی، روشن بود که حضور در اینترنت برای برنده شدن حیاتی است و این امر تمرکز بر برخط بودن سیاسی را افزایش داد. آرشیو اینترنت با کتابخانه کنگره برای گردآوری سایت‌های سیاسی شریک شد.

این اقدام، اولین پروژه آرشیو با کتابخانه کنگره بود و برای بسیاری از کارکنان آن به معنی انتقال از

Streamable -1

International Internet Preservation Consortium (IIPC) <http://netpreserve.org> -2

Rick Prelinger -3

یک پروژه تجربی به یک مؤسسه ثابت بود.

فکر ایجاد دسترسی به آثار ناپایدار نگهداری شده باعث حرکت و جنبش آئی و آرشیو تصاویر متحرک (1) در سال 2000 تأسیس شد. در حال حاضر، این آرشیو، با پروانه Creative Commons، شامل فیلم هایی از آرشیو Prelinger مجموعه ای بیش از 1900 فیلم ناپایدار (2) (پیام های بازرگانی، آموزشی، صنعتی و آماتور) است. این مجموعه، در حال حاضر دارای بیش از 10 درصد تمام فیلم های ناپایدار تولید شده در سالهای 1927 و 1987 در آمریکاست و یکی از کامل ترین و متنوع ترین مجموعه فیلم های ژانرهای است که تعداد زیادی از آن ها نگهداری نشده اند.

2001: دسترسی از طریق ماشین Wayback: آرشیو یازده سپتامبر

طی یک بازه زمانی یکساله از مارس 2000 تا مارس 2001، آرشیو، اندازه موجودی خود را 3 برابر کرد و به بیش از 40 ترابایت رساند. در این دوره، آرشیو هر ماه تقریباً 10 ترابایت رشد می کرد.

عکس

یک پروژه تجربی به یک مؤسسه ثابت بود.

فکر ایجاد دسترسی به آثار ناپایدار نگهداری شده باعث حرکت و جنبش آئی و آرشیو تصاویر متحرک^۱ در سال ۲۰۰۰ تأسیس شد. در حال حاضر، این آرشیو، با پروانه Creative Commons، شامل فیلم‌هایی از آرشیو Prelinger، مجموعه‌ای بیش از ۱۹۰۰ فیلم ناپایدار^۲ (پیام‌های بازرگانی، آموزشی، صنعتی، و آماتور) است. این مجموعه، در حال حاضر، دارای بیش از ۱۰ درصد تمام فیلم‌های ناپایدار تولید شده در سال‌های ۱۹۲۷ و ۱۹۸۷ در آمریکاست و یکی از کامل‌ترین و متنوع‌ترین مجموعه فیلم‌های ژانرهای است که تعداد زیادی از آنها نگهداری نشده‌اند.

۲۰۰۱: دسترسی از طریق ماشین Wayback: آرشیو یازده سپتامبر

طی یک بازه زمانی یکساله از مارس ۲۰۰۰ تا مارس ۲۰۰۱، آرشیو، اندازه موجودی خود را ۳ برابر کرد و به بیش از ۴۰ ترابایت رساند. در این دوره، آرشیو، هر ماه تقریباً ۱۰ ترابایت رشد می‌کرد.

جدول ۲. مجموعه‌های آرشیو در مارس ۲۰۰۱

اندازه	واحد	مجموعه
۴۰ ترابایت	۴ بیلیون صفحه	۱۹۹۶- ماه ۳ سال ۲۰۰۱
۲ ترابایت	۲۰۰ میلیون صفحه	آرشیو انتخابات ۲۰۰۰
۰/۵ ترابایت	۱۶ میلیون پستینگ	۱۹۹۶-۱۹۹۸، ۲۰۰۰- ماه ۳ سال ۲۰۰۱
۰/۵ ترابایت	۳۶۰ فیلم	فیلم‌های آرشیو: حدود ۱۹۰۳- حدود ۱۹۷۳
<۰/۱ ترابایت	۵۰۰۰ صفحه	Arpanet : اسناد تاریخی

اما سال ۲۰۰۱ برای بسیاری از سازمان‌های دارای فناوری بالا در محدوده سانفرانسیسکو، سال سختی بود. سقوط بازار سرمایه، واگذاری صدها شرکت محلی، و حمله به مرکز تجارت جهانی در نیویورک بر کارهای آرشیو اثر گذاشت. به‌ویژه، از دست رفتن مشاغل «های‌تک»^۳ در واقعه ۱۱ سپتامبر باعث تمرکز بر آن شد. همانگونه که برای توان آرشیو نیز در تهیه تصاویر متحرک و پاسخ به این وقایع حکم نوعی آزمون را داشت. در اوایل ۲۰۰۱، شاید مهم‌ترین پرسش پیش روی آرشیو این بود که چگونه به بهترین وجه، دسترسی به مجموعه را فراهم کند. داده‌های بسیاری، به‌طور مستقیم، از طریق خدمات الکسا در اختیار عموم قرار می‌گرفت، ولی دسترسی مستقیم به مجموعه‌ها هنوز نیازمند مهارت‌های برنامه‌ریزی یونیکس^۴ بود.

برنامه‌ریزان الکسا، در قراردادی با آرشیو، برنامه ماشین Wayback را ساختند که خدمات دسترسی

1. Moving Images Archive
2. ephemeral
3. high tech
4. Unix

جدول ۲. مجموعه‌های آرشیو در مارس ۲۰۰۱

اما سال ۲۰۰۱ برای بسیاری از سازمان‌های دارای فناوری بالا در محدوده سانفرانسیسکو، سال سختی بود سقوط بازار سرمایه، واگذاری صدها شرکت محلی، و حمله به مرکز تجارت جهانی در نیویورک بر کارهای آرشیو اثر گذاشت. به‌ویژه، از دست رفتن مشاغل «های‌تک»⁽³⁾ در واقعه ۱۱ سپتامبر باعث تمرکز بر آن شد. همانگونه که برای توان آرشیو نیز در تهیه تصاویر متحرک و پاسخ به این وقایع حکم نوعی آزمون را داشت. در اوایل ۲۰۰۱، شاید مهم‌ترین پرسش پیش روی آرشیو این بود که چگونه به بهترین وجه، دسترسی به مجموعه را فراهم کند. داده‌های بسیاری، به‌طور مستقیم، از طریق خدمات الکسا در اختیار عموم قرار می‌گرفت، ولی دسترسی مستقیم به مجموعه

ها هنوز نیازمند مهارت های برنامه ریزی یونیکس (4) بود.

برنامه ریزان الکسا، در قراردادی با آرشیو ، برنامه ماشین Wayback را ساختند که خدمات دسترسی

ص: 69

Moving Images Archive -1

ephemeral -2

high tech -3

Unix -4

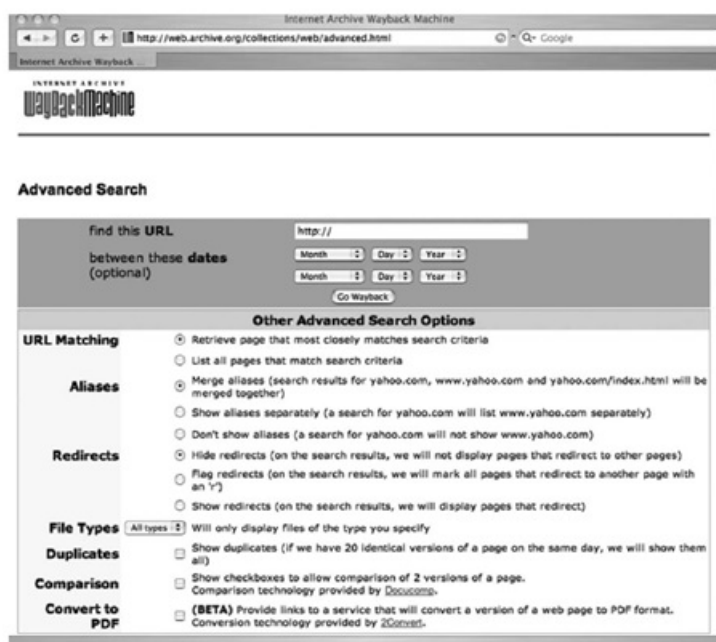
به محتوای آرشیو را بر پایه URL ها می داد 24 اکتبر 2001، ماشین Wayback به کار افتاد و دسترسی به بیش از 10 بلیون صفحه وبی آرشیو شده و 100 ترابایت داده ممکن شد.

در آن زمان، داده ها روی Hewlett Packar ذخیره می شد و سرورهای uslab.com از سیستم های عامل Linux و FreeBSD استفاده می کرد. هر رایانه، حدود 512 مگابایت حافظه داشت و به طور کلی بیش از 300 گیگابایت روی دیسک های IDE بود.

عکس

۷۰ مدیریت منابع اطلاعاتی وب

به محتوای آرشیو را بر پایه URL ها می داد. ۲۴ اکتبر ۲۰۰۱، ماشین Wayback به کار افتاد و دسترسی به بیش از ۱۰ بلیون صفحه وبی آرشیو شده و ۱۰۰ ترابایت داده ممکن شد.
در آن زمان، داده ها، روی Hewlett Packar ذخیره می شد و سرورهای uslab.com از سیستم های عامل Linux و FreeBSD استفاده می کرد. هر رایانه، حدود ۵۱۲ مگابایت حافظه داشت و به طور کلی بیش از ۳۰۰ گیگابایت روی دیسک های IDE بود.



تصویر ۹.۱. گزینه های فعلی جستجو برای Wayback machine

پروژه مهم دیگر سال ۲۰۰۱، آرشیو ۱۱ سپتامبر^۱ بود. با همکاری کتابخانه کنگره، آرشیو، تصاویری از بیش از ۳۰۰,۰۰۰ وبگاه منتخب را از ۱۱ سپتامبر ۲۰۰۱ تا ۱ دسامبر ۲۰۰۱ و صدها ساعت پخش اخبار را گردآوری کرد.

۲۰۰۲: کتابخانه اسکندریه^۲، کتابخانه سیار، و حق مؤلف

در سال ۲۰۰۲، آرشیو، ۵ برنامه مهم دیگر یعنی افزایش مجموعه ها، گزینه های دستیابی به این مجموعه ها، همکاری با سایر سازمان ها و تعیین راهکارها را بر عهده داشت.

1. September 11 Archive
2. The library of Alexandria

تصویر 1.9. گزینه های فعلی جستجو برای Wayback machine

پروژه مهم دیگر سال 2001، آرشیو 11 سپتامبر (1) بود. با همکاری کتابخانه کنگره، آرشیو، تصاویری از بیش از 30,000 وبگاه منتخب را از 11 سپتامبر 2001 تا 1 دسامبر 2001 و صدها ساعت پخش اخبار را گردآوری کرد.

2002: کتابخانه اسکندریه، کتابخانه سیار، و حق مؤلف

*کتابخانه اسکندریه، کتابخانه سیار، و حق مؤلف (2)

در سال 2002، آرشیو، 5 برنامه مهم دیگر یعنی افزایش مجموعه ها، گزینه های دستیابی به این مجموعه ها، همکاری با سایر سازمان ها و تعیین راهکارها را بر عهده داشت.

ص: 70

September 11 Archive -1

The library of Alexandria -2

نخستین و بزرگ ترین پروژه ایجاد سایت قرینه در کتابخانه اسکندریه در مصر بود. سرورها و بیش از 100 ترابایت داده، که بیش از 5 میلیون دلار ارزش داشتند به مصر برده شد و برای افتتاح کتابخانه در ماه آوریل نصب شد.

دومین پروژه مهم، ایجاد کتابخانه سیار اینترنتی (1) بود که برای نمایش چگونگی ترکیب اسکن های الکترونیکی کتاب ها، چاپ بر اساس فناوری مورد تقاضا، و اینکه ارتباط شبکه ای ماهواره ای چگونه می تواند مناسب کتابخانه ای هزار جلدی در پشت یک ون باشد، طراحی شد پروژه یک میلیون کتاب (2) Million Books Project، در تابستان 2002، با شراکت کارنگی ملون شروع شد. هدف این بود که حداقل یک میلیون کتاب دیجیتالی شود و به صورت رایگان روی اینترنت قرار گیرد. با تشویق سوزان مبارک (3)، آرشیو، ساخت یک کتابخانه سیار را در ایالات متحده شروع کرد و با همکاری دیگران، در هند و کنیا، نیز نمونه هایی از آن ایجاد شد.

سومین پروژه مهم مربوط به سیاست گذاری بود. 30 سپتامبر 2002، در تلاشی برای افزایش آگاهی عمومی درباره اهمیت موارد راهکاری حق مؤلف و کتابخانه سیار اینترنتی، آرشیو به سفری در سر تاسر کشور و چاپ و توزیع کتاب های رایگان مبادرت کرد کتابخانه سیار را در حیاط ساختمان دادگاه عالی ایالات متحده پارک کرد و به چاپ کتاب پرداخت؛ جایی که قضات، در 9 اکتبر مباحثه میان الدرد (4) و اشکرافت (5) را شنیدند این حادثه مهمی بود که بر اساس آن تصمیم گرفته شد چه تعداد کتاب باید بخشی از کتابخانه دیجیتالی کتابخانه سیار و سایر کتابخانه های دیجیتالی در ایالات متحده باشد. متأسفانه، الدرد شکست خورد و تصمیم حق مؤلف مؤثر واقع افتاد؛ ولی پروژه کتابخانه سیار شکوفا شد و در نهایت، به صورت کتابخانه ای غیر انتفاعی درآمد. اخیراً، دولت هند، ساخت 25 کتابخانه سیار برای استفاده در سر تاسر کشور هند را شروع کرده است.

چهارمین حوزه فعالیت شامل ایجاد اولین آرشیو مجموعه کتاب و موسیقی است در ژوئن اولین مجموعه های کتاب، به صورت پیوسته در اختیار گذاشته شد در آگوست آرشیو موسیقی زنده، شامل مجموعه ای از اجراهای کنسرت - که به طور قانونی قابل بارگیری بودند - به صورت پیوسته درآمد.

پنجمین پروژه مهم، تأسیس کتابخانه دیجیتالی بین المللی کودکان (6)، با شرکت دانشگاه مری لند (7)، بود که از سوی کتابخانه کنگره، NSF، IMLS بنیاد کال / اوستین (8)، شرکت سیستم های ادوب (9)، بنیاد مرکل (10)، و اکتاوو (11) پشتیبانی شد. ICDL، بر ماهیت ذاتی اینترنت، مبتنی بر فراهم آوری دسترسی مستقیم و جهانی به محتوای کیفی برای کودکان متمرکز بوده و هست.

ص: 71

Internet Bookmobile -1

MBP)، -2

Suzanne Mubarak -3

Eldred -4

Ashcroft -5

International Children's Digital Library -6

University of Maryland -7

Kahle/Austin Foundation -8

.Adobe System Inc -9

Markle Foundation -10

Octavo -11

در پایان سال 2002، آرشیو، برای مراقبت از یکپارچگی آرشیوهای دیجیتالی از طریق استاندارد سازی معیارها، و جلوگیری از جابه جا یا غیر قابل دسترس شدن مواد تلاش کرد. در نشست در یوسی برکلی، نمایندگان از آرشیو با سایر کتابداران دیجیتالی، به منظور تکمیل Okland Archive Policy ملاقات کردند که روند انتقال مواد را بر اساس قانون یا بر اساس خواست مالکان سایت و سایرین با جزئیات بیان می کرد.

2003: گسترش دستیابی ما به کتابخانه های ملی و مؤسسات آموزشی

در سال 2003، آرشیو، رسیدن به کتابخانه های ملی و مؤسسات آموزشی در سرتاسر دنیا را ادامه داد. آرشیو، با همکاری کنسرسیوم حفاظت بین المللی اینترنتی (IIPC) از نزدیک با سازمان های شریک روی استانداردهای جدید و یک خزشگر جدید منبع باز شروع به کار کردن کرد.

در جولای 2003، آرشیو، به کنسرسیوم حفاظت بین المللی اینترنتی برای آغاز به کار کمک کرد. گروهی متشکل از 12 کتابخانه ملی برای توسعه استانداردها، ابزار و راهکارهایی برای تهیه حفاظت و دسترس پذیری دانش و اطلاعات از اینترنت برای نسل های آینده در همه جا، ارتقای تبادل جهانی و ارتباطات بین المللی تلاش می کردند. برای انجام این رسالت، IIPC برای رسیدن به اهداف زیر کار می کند:

- رسیدن به مجموعه ای غنی از محتوای اینترنتی از سرتاسر دنیا که به گونه ای حفاظت شوند که بتوانند آرشیو و محافظت شده و در لحظه قابل دسترسی باشند؛

- توسعه و استفاده از ابزار عادی فناوری ها و استانداردهایی که ایجاد آرشیوهای اینترنتی را ممکن می سازند؛

- تشویق و حمایت از کتابخانه های ملی در همه جا برای آرشیو اینترنتی و حفاظت.

IPC در کتابخانه ملی فرانسه با 12 مؤسسه همکار مجاز شروع به کار کرد. اعضا موافقت کردند به اتفاق هزینه ها را پرداخت کرده و در برنامه ها و کارگروه ها برای رسیدن به اهداف مذکور شرکت کنند. موافقت اولیه برای 3 سال بود. طی این پروژه اعضاء محدود به مؤسسات مجاز بودند.

در سال 2003، آرشیو بودجه قابل توجهی از خارج یعنی از سازمان های دیگر از جمله بنیادهای هیولیت واسلون (1) دریافت کرد و شروع به کار بر مجموعه های خاص کرد. کاهش های بعدی هزینه ذخیره روی دیسک و پهنای باند اینترنت آرشیو را به عرضه دائم «پهنای باند نامحدود، برای همیشه و رایگان» برای سازمان ها و افراد با مواد دیجیتالی، راهبر شد.

این امر به شراکت با Etree منجر شد. سازمانی که داوطلبانه در سال 1998 ایجاد شد تا تجارت آزاد و قانونی کنسرت های موسیقی زنده را امکان پذیر سازد. حاصل همکاری با Etree این شد که آرشیو اکنون میزبان بیش از 15000 کنسرت موسیقی زنده است.

آرشیو، برای حمایت از نیازهای رو به رشد هم ذخیره و هم پهنای باند، یک مرکز جدید در سانفرانسیسکو باز کرد. این مرکز داده جدید، از طریق پیوند 1 Gbps به اینترنت متصل است و بیش از 1500 رایانه شخصی را که از لینوکس استفاده می کنند، میزبانی می کند.

اشاره

*آرشیو اروپا و پتاباکس (1) (2)

آرشیو، در سال 2004 شروع به انتقال داده ها به سومین نسل سخت افزاری خود، موسوم به پتاباکس، کرد. طرح پتاباکس، بر پایه سخت افزاری rack-mounted و سیستم عامل لینوکس، ذخیره RAID برای هر ترابایت به مبلغ تقریبی 2000 دلار یا 2 میلیون دلار برای هر پتابایت پیشنهاد کرد.

نخستین نصب این طرح جدید در آمستردام در آرشیو تازه تشکیل شده اروپا بود، مؤسسه ای که قرار بود پاسخگوی نیازهای جامعه اروپا باشد و آرشیو اینترنت با سایر شرکای اروپایی در اولین سال تأسیس، آن را حمایت می کرد. نصب آن در آمستردام برای ایجاد قرینه ای برای مجموعه های اسکندریه و سانفرانسیسکو است. ایجاد شبکه ای از مؤسسات مستقل در سرتاسر دنیا که هر یک قادر است به طور مستقل عمل کند به پیشگیری از نابودی های فاجعه آمیز اطلاعات کمک خواهد کرد.

همچنین، در سال 2004، کنسرسیوم حفاظت بین المللی اینترنتی (Heritrix 3)، یعنی خزشگر وی منبع باز، قابل توسعه با قابلیت تغییر اندازه وی، با کیفیت آرشیوی و بر پایه جاوا را شروع کرد.

در مسیر توسعه مجموعه در سال 2004، از طریق استخدام کارکنان بیشتر، تکمیل پروژه های اسکن کتاب و فیلم و دریافت داده از سایر مؤسسه ها، گام های مهمی رو به جلو برداشته شد.

آینده

پیشرفت در رایانه و ارتباطات این امکان را فراهم ساخته که ذخیره کتاب، صفحه موسیقی، فیلم، بسته های نرم افزاری و تمام صفحات عمومی وب که تاکنون ساخته شده اند، هزینه - سودمندی داشته باشند و دسترسی به این مجموعه ها از طریق اینترنت برای همه از جوان تا پیر در سرتاسر دنیا فراهم شود.

همانگونه که در اعلامیه حقوق بشر، ماده 19 آمده است: «هر انسانی حق آزادی بیان و عقیده را دارد. این حق شامل آزادی داشتن باور و عقیده بدون نگرانی از مداخله و مزاحمت و حق جست و جو، دریافت و انتشار اطلاعات و نظرها از طریق هر رسانه ای بدون ملاحظات مرزی است».

برای سال های آینده، رسالت آرشیو مشخص است: ایجاد کتابخانه ای جهانی که تمام دانش را به سهولت در اختیار هر مرد و زن و کودک در سرتاسر جهان قرار می دهد.

منابع

Gray, J. Shenoy, P. (1999). Rules of Thumb in Data Engineering. Microsoft

Technical Report, MS-TR-99-100

Masanès, J. (2006). Web archiving: issues and methods. In J. Masanès (Ed.), Web Archiving. Springer, Berlin Heidelberg New York

Ubois, J. (2002). The Oakland Archive Policy. Recommendations for Managing Removal Requests and Preserving Archival Integrity

ص: 73

European Archive -1

Petabox -2

Heretrix -3

پیشرفت در رایانه و ارتباطات این امکان را فراهم ساخته که ذخیره کتاب، صفحه موسیقی، فیلم، بسته های نرم افزاری و تمام صفحه های عمومی وب که تاکنون ساخته شده اند، هزینه - سودمندی داشته باشند و دسترسی به این مجموعه ها از طریق اینترنت برای همه، از جوان تا پیر، در سر تاسر دنیا فراهم شود. برای سال های آینده، رسالت آرشیو مشخص است: ایجاد کتابخانه ای جهانی که تمام دانش را به سهولت در اختیار هر مرد و زن و کودک در سر تاسر جهان قرار دهد. در مقاله حاضر بحث آرشیو اینترنت و آرشیو در محیط اینترنت، دسترس پذیری دراز مدت تمام دانش برای همه افراد جهان پیشرفت های فناورانه ایجاد آرشیو اینترنتی مورد بررسی قرار می گیرد. سپس برنامه ها و پروژه های نمونه آرشیو اینترنتی در دنیا آرشیو اروپا و پتاباکس معرفی می شوند.

*کاربرد وب و مطالعات مربوط به آن (1)

استیو جونز (2) | گمیل جانسون (3) (دانشگاه ایلینویز شیکاگو)

ترجمه: دکتر سید مهدی طاهری (4) | سید محمد موسوی (5)

خلاصه

اشاره

در سال 2002، مرکز کتابخانه رایانه ای پیوسته (OCLC) برآورد کرد که بیش از سه میلیون وبگاه در شبکه وب جهان گستر در دسترس عموم قرار دارد (اونیل و همکاران 2003، پاراگراف 9). در دنیای مادی اطلاعات، این تعداد وبگاه نزدیک به 14 تا 28 میلیون خواهد بود که با حجم کتاب های موجود در بزرگ ترین کتابخانه های دنیا برابری یا حتی از تعداد آن تجاوز می کند. این مجموعه گسترده اطلاعات، مجموعه ای بی پایان و دسترس پذیر از داده های دیداری و شنیداری را در اختیار پژوهشگران علاقه مند به مطالعه در حوزه فعالیت های پیوسته قرار می دهد. حجم خالص منابع دسترس پذیر بر روی وب همواره چالش هایی را به وجود آورده است چه چیزی را انتخاب و چگونه مطالعه کنیم؛ با این حال، شبکه وب ثابت کرده است این حجم اطلاعات برای کاربران دلهره آورتر هم شده است. در این مقاله، برخی رویکردهای روش شناختی را که پژوهشگران برای مطالعه در حوزه وب به کار گرفته اند، مرور خواهیم کرد.

ص: 75

Web Use and Web Studies: in Masanes, Julien (ed.), Web Archiving. Berlin Heidelberg New York: - 1
.Springer. pp.55-67

Steve Jones -2

Camille Johnson -3

4- عضو هیئت علمی پژوهشگاه علوم و فرهنگ اسلامی

5- کارشناس ارشد کتابداری و اطلاع رسانی پژوهشگاه علوم و فرهنگ اسلامی

هدف از این کار، طبقه‌بندی جامع روش‌شناسی‌ها نیست به جای آن امیدواریم با درک روش‌های به کار گرفته شده در مطالعه وب و مطالعه کاربرد آن، افرادی که درصدد آرشو کردن و نگهداری وب هستند بتوانند نیازهای جامعه علمی را هرچه بهتر درک کنند.

شبکه وب عبارت است از طیف گسترده‌ای از انواع مواد که تنوع آن‌ها را با دو بُعد می‌توان بهتر درک کرد. نخست اینکه خود وب یک رسانه است نه محتوا برای روشن تر شدن مطلب باید گفت که وب هم رسانه‌ای است که محتوا را از طریق پروتکل‌های متعدد (نظیر HTTP) منتقل می‌کند و هم «ظرفی» برای محتواست که محتوا را شکل داده و به کاربران خود ارائه می‌دهد. هر چند ارائه محتوا به تشخیص کاربر وابسته است و گذشته از آن توسط ابزارهای دیداری مورد استفاده (مرورگرها و برنامه‌های کاربردی دیگر) شکل می‌گیرد. به عبارت دیگر، اگر چه محتوا ممکن است صفحه وب یکسانی باشد، دو کاربری که از مرورگرهای متفاوت، یا از تنظیمات متفاوتی از یک مرورگر استفاده می‌کنند، ممکن است در نهایت صفحه‌های متفاوتی را ببینند. دیگر اینکه، برخلاف دنیای مواد آنالوگ، همان تعریف رسانه برای ذخیره‌سازی وب مورد تردید است. وبگاه‌ها را می‌توان به صورت محلی ذخیره یا پنهان کرد، اطلاعات آن‌ها را منعکس نمود، یا ممکن است مجازی یا موقتی باشند همان‌طور که این ویژگی‌ها برای وب کم‌ها نیز صادق است.

به طور کلی، پژوهش درباره اینترنت و وب، توسط پژوهشگرانی با تخصص‌های مختلف از جمله متخصصان رشته‌های زبان‌شناسی، روزنامه‌نگاری، علوم سیاسی، مدیریت بازرگانی، جغرافیا، تبلیغات، ارتباطات، هنر، و دیگر رشته‌ها انجام می‌شود. از آن‌جا که با طیف وسیعی از سنت‌های آموزشی سر و کار داریم، انواع موادی که پژوهشگران از آن‌ها بهره می‌گیرند بسیار متنوع است و به همان اندازه نیز در روش‌های پژوهش تنوع مشاهده می‌شود. آن‌چه به عنوان تحلیل «قوم‌نگاری» متون وب، برای پژوهشگر بازاریابی به حساب می‌آید، ممکن است از دیدگاه پژوهشگر رشته ارتباطات تحلیل محتوای کیفی تعبیر گردد. بنابراین هدف دیگر این مقاله نه تنها ارائه مقوله‌های ثابت محتوای وب یا طرح‌هایی برای رویکردهای روش‌شناختی است - آن‌گونه که برای مطالعات وب به کار می‌روند - بلکه ارائه طیفی از تعابیر و کاربردهای این روش‌ها نیز مورد نظر است، به نحوی که مفیدترین روش‌ها را برای پژوهشگران حوزه وب به اثبات رساند.

1- تحلیل محتوا

تحلیل محتوای یکی از رایج‌ترین روش‌های مطالعه بر روی حوزه وب است. پژوهشگری که از روش تحلیل محتوا استفاده می‌کند، محتوای وب را - هم محتوای متنی و هم تصاویر - بر اساس معیارهای خاص رمزگذاری کرده و آن‌ها را درون مقوله‌ها یا موضوع‌های مرتبط قرار می‌دهد؛ به عبارت دیگر، این نوع تحلیل بیشتر بررسی محتوای وب است تا کاربران وب. تحلیل محتوا در میان مطالعات وب، به عنوان ابزار مقایسه مورد استفاده قرار گرفته است به نحوی که این ابزار به پژوهشگر اجازه می‌دهد تا مقایسه‌های معنی‌داری از محتوا میان متون مشابه وب انجام دهد. در بررسی وبگاه‌های سازمان‌های ضد جهانی‌سازی

از روش تحلیل محتوا استفاده شد تا نشان داده شود که آیا میان این سایت ها انسجام پیام و مقصود وجود دارد (ون الست و، والگریو 2002)؟ محتوا، بر اساس کارکرد اصلی خود در چهار حوزه رمزگذاری و مرتب شده است ارائه اطلاعات درباره سازمان مسائل ضد جهانی سازی افزایش تعامل با مؤسسان گروه و اعضای دیگر و افزایش بسیج افراد در مورد مسائلی چون اهدای پول که در یک فراخوان اینترنتی اعلام شده است و الست و والگریو، از طریق تحلیل خود به این نتیجه رسیدند که در حقیقت سازمان های ضد جهانی سازی پیوسته، عموماً به وسیله وبگاه های خود و از روش های مشابهی برای آگاه نمودن و مشغول کردن اعضای خود استفاده می کردند.

مطالعه تطبیقی دیگری به تحلیل محتوای وبگاه های مربوط به ایستگاه های رادیویی پرداخت، تا چگونگی پاسخگویی صنعت تجاری رادیو را به اندازه دسترس پذیری وب جهان گستر برای ارتقای ایستگاه های خود مشخص کند. این ارزیابی با تجزیه و تحلیل انواع اطلاعات کاربر - محور، که در وبگاه های ایستگاه هایی رادیویی فراهم شده بود برای مثال نقشه های ترافیکی و لاگ های برنامه ها، استفاده از وبگاه ها به عنوان ابزارهایی برای بهبود ایستگاه خود (مانند اطلاعات رقابتی و بایوس های DG)، و ترکیب ویژگی های تعاملی آن ها (نظیر نشانی ایمیل ها برای اعضای ایستگاه ها و نظرسنجی های شنوندگان) به دست آمد (پیتز و هارمز 2003). این مطالعه در حوزه وب منبع اطلاعاتی ارزشمندی را برای ارزیابی کاربردهای کنونی و بالقوه وب در اختیار ایستگاه های رادیویی قرار داد.

همچنین، روش تحلیل محتوا، برای مطالعه اثرات تغییرات خط مشی های سازمانی بر وبگاه های مرتبط به کار گرفته شده است. برای مثال، خط مشی های آموزشی، اهمیت الحاق فناوری اطلاعات و ارتباطات (ICT) به مأموریت های سازمانی مدارس راهنمایی را افزایش دادند، مطالعه ای، پیشرفت های صورت گرفته برای نیل به این هدف را همخوان با طرح ریزی انجام شده توسط شبکه ملی یادگیری در بریتانیا از طریق تحلیل محتوا در وبگاه های 150 مدرسه راهنمایی مورد ارزیابی قرار داد (هسکت و سولیوان 1999). تصاویر و متون رمزگذاری شدند، و به شناسایی وبگاه ها در یکی از مقوله های روشی زیر پرداختند: فعال، دانش آموز-محور، و منفعل. پژوهشگران از طریق این مقوله ها و تحلیل پرونده های پژوهشی مدارس، رابطه ای را بین تعهد به یکپارچگی فناوری اطلاعات و ارتباطات (ICT) و سرمایه های سازمانی درک نمودند؛ هر قدر سرمایه اقتصادی و اجتماعی مدرسه بیشتر باشد احتمال اینکه محتوای وبگاه آن منعکس کننده گرایش های مثبت و فعال الحاق فناوری اطلاعات و ارتباطات با مأموریت های مدرسه باشد بیشتر خواهد بود.

کمیسیون بازرگانی فدرال امریکا، از روش تحلیل محتوا، برای ارزیابی پیاده سازی برنامه های حفظ امنیت اطلاعات و حریم خصوصی پیوسته بر روی وبگاه های بازرگانی استفاده کرده است (Milne and Culnan). مطالعه انجام شده از سوی این کمیسیون، دارای تحلیل دراز مدت چهار نظر سنجی وب با گستره زمانی 1998 تا 2001 بود از معیارهایی نظیر حفظ یا حذف اعلان ها، با توجه به افشای اطلاعات بازدیدکنندگان با اشخاص ثالث و استفاده از «کوکی ها»، مجموعه ای از اطلاعات کاربران به غیر از نشانی ایمیل ها، و به کارگیری امنیت اطلاعاتی در سایت ها مورد استفاده قرار گرفت تا مشخص شود که آیا

وبگاه ها ملزومات شیوه های استفاده بی طرفانه از اطلاعات را تأمین می نمایند؟ مقایسه نتایج به دست آمده از تحلیل انجام شده در سال 2001 با تحلیل های سال های پیشین نشان می دهد، در میان وبگاه های تجاری که به گردآوری اطلاعات شخصی کاربران خود دست می زنند، روند افزایشی در اجرای شیوه های استفاده بی طرفانه از اطلاعات وجود داشته است.

2- بررسی ها

بررسی های مربوط به وب را می توان به دو روش متمایز تعریف کرد: بررسی هایی که بر روی وب انجام می گیرد، و بررسی های منتشر شده در وب. بررسی های نوع نخست که از تداول بیشتری برخوردار است به وسیله پژوهشگرانی انجام شده است که مطالعه درباره وبگاه ها را به منظور دسترس پذیری برای جمعیت علاقه مند ارسال می کنند. در بیشتر موارد، هدف از این نوع بررسی ها، گردآوری اطلاعاتی درباره موارد استفاده کاربران از اینترنت (1) است (مانند بررسی کاربران وب که در مرکز گرافیک، دیداری سازی، و کاربرد پذیری (The [Graphic, Visualization, and Usability Center](#) (2) انجام گرفته است (1998). GVU، از سال 1994، مطالعات مداومی درباره کاربرد اینترنت انجام داده است. این مرکز، شرکت کنندگانی را از طریق ثبت گروه های خبری مرتبط با اینترنت، تابلوهای تبلیغاتی در رسانه های خبری، وبگاه های موتورهای کاوش و اطلاعیه های موجود در رسانه های ناپیوسته نظیر مجلات و روزنامه ها به کار می گیرد. چندین بررسی متمرکز توسط GVU از سال 1994 تولید شد که از آن میان می توان به بررسی سرشماری عمومی، بررسی مرتبط با سرشماری فناوری- که اطلاعاتی را نظیر سرعت ارتباط کاربران و مرورگر مورد استفاده گردآوری کرد، بررسی مربوط به گرایش کاربران به سوی حریم خصوصی و امنیت اینترنتی، بررسی عمومی رایانه ها، کاربرد اینترنت و وب، و نیز بررسی جست و جوی محصولات و فعالیت های مربوط به خرید به صورت پیوسته اشاره کرد. جامعه آماری، معمولاً بالای 1000 [نفر] بودند که بررسی مربوط به سرشماری عمومی مشتمل بر 5000 شرکت کننده بود.

سازمان ها نیز بررسی های مربوط به وب را ابزار مفیدی برای ارزیابی کیفیت تجربیات کاربرانی می دانستند که از وبگاه ها یا کارکرد های وب بهره می برند. چنین مطالعه ای، بررسی بازدیدکنندگان حدود 450 وبگاه موزه ها را هدایت می کند به نحوی که هم اطلاعات مربوط به سرشماری، و هم پاسخ های مربوط به کیفیت و قابلیت استفاده از این وبگاه ها را گردآوری می نماید (Sarraf 1999). بسیاری از موزه ها به تدریج وب گاه هایی را تهیه کردند تا ابزاری برای ارائه اطلاعات درباره مجموعه های موجود در موزه ها و نیز ایجاد ارتباطی تعاملی و قابل دسترس برای بازدیدکنندگان کنونی و پیشین باشد.

ص: 78

1- با دسترسی 59 درصد از افراد بالغ آمریکایی به اینترنت در سال 2002، پژوهشگران در آمریکا نیز بررسی های مبتنی بر وب را به جای انواع بررسی های سنتی از قبیل بررسی های تلفنی و کاغذی درباره موضوعاتی به جز کاربردهای اینترنت و نگرش های موجود مد نظر قرار دادند به نظر می رسد پژوهشگران علاقه مند به گردآوری هنگام واکنش ها نسبت به وقایع اخیر به ویژه به بررسی های اینترنتی هستند، به نحوی که این علاقه مندی در بررسی های انجام گرفته درباره اثرات حملات تروریستی نیویورک در 11 سپتامبر سال 2001 بر شرکت کنندگان مشهود است (برای مثال لی و همکاران 2003) درست 9 روز پس از این حملات، گروهی از پژوهشگران توانستند بررسی از طریق ارزیابی وب-پایه پاسخ های روان شناختی آمریکایی ها نسبت به این واقعه تلخ را توزیع کنند.

2-، GVU)

به منظور ارائه بازخوردی جامع برای موزه های مورد بحث به پاسخ های کیفی و کمی نیاز بود. کارزول و ونکاتش (2002)، با استفاده از بررسی وب محور توانستند پاسخ های ارزیابی بیش از 500 دانشجوی مقطع کارشناسی را که در کلاس های پیوسته غیر همزمان شرکت داشتند درخواست کنند. آن ها حمایتی نسبی را در مورد فرضیه های خود یافتند مبنی بر اینکه پذیرش، و در آینده قصد استفاده از فناوری در دوره ای پیوسته و غیر همزمان می تواند به صورت مثبت متأثر از گرایش ها و درک دانشجویان از فناوری باشد. در تمامی موارد از جمله GVU، مزایای استفاده از بررسی مبتنی بر وب آشکار بود: جمعیت های علاقه مند در بررسی ها همان کاربران وب بودند، و داده های گردآوری شده با فعالیت های وب ارتباط مستقیم داشت.

3- تحلیل بلاغی

*تحلیل بلاغی (1)

زیر مجموعه دیگری از پژوهش در حوزه وب نگاهی جدی به متون وب دارد تا راهکارهای متقاعد کننده ای را از طریق تحلیل بلاغی مشخص کند. همان طور که وارنیک توصیف می کند:

روش انتقادی بلاغی این مسئله را در نظر دارد که چگونه متن به برخی عناصر موجودیت می بخشد- در حالی که با برخی دیگر مخالفت می کند، چگونه ساختارهای گزارشی تجربه مرورگر را به روش های مشخص پیکربندی می کنند، و اینکه چگونه گفتمان، تمایلات و عادات ذهنی مخاطب خود را نقش می دهد.

وی، در بررسی خود به وبگاه های سیاسی در جریان مبارزات انتخاباتی رئیس جمهوری آمریکا در سال 1996 به ویژه راهکارهای بلاغی وبگاه های مقلد پرداخت - وبگاه هایی که از طراحی و محتوای وبگاه های قانونی مربوط به مبارزات انتخاباتی تقلید می کردند تا در صورت امکان از اعتبار آن سایت ها بکاهند و تا حدی آن ها را مورد تمسخر قرار دهند. بر اساس نتیجه گیری، وارنیک چنین سایت هایی بر گزارش های نادرست در دولت، تحریک مشارکت سیاسی از طریق ویژگی های تعاملی - نظیر شکایاتی که هرگز به حزب مورد نظرشان تسلیم نشد و اشاراتی درباره سوابق سیاسی و جنایی - که نامزدها بر آن ها متکی هستند تا خوانندگان شان را به پذیرش بینش بلاغی خود مجاب کنند - تکیه می کنند. در نهایت، وارنیک، این راهکارها را ریاکارانه می داند؛ زیرا از رفتارهای غیر اخلاقی مشابه بسیاری که ظاهراً از سوی نامزدهای مورد نظر به کار گرفته می شود بهره می گیرند.

از تحلیل بلاغی وبگاه هایی که برای شناسایی اجتماعات بلاغی چنین استفاده می گردد که گروه هایی از وب برای شکل گیری جهان بینی خود استفاده می کنند که ممکن است همسو با جریانات اخیر باشد یا نباشد کروبر، به بحث و تبادل نظر در این مسئله می پردازد که گروهی از وب گاه های شکل گرفته توسط مادرانی با گرایش های فمینیستی، خط مشی [خود] را به مخاطبان القا می کند تا به طرز مؤثری در برابر تصورات منفی فمینیستی در رابطه با مادر بودن مقاومت ورزند. زنان، توانستند از طریق متون و تصاویر موجود در وب گاه های شان دوباره به بحث و تبادل نظر در خصوص فمینیسم پردازند و برای درک معنای مادری به عنوان بخش پر نفوذ و حتی لازم فمینیست بودن به بحث و جدل می پردازند.

بررسی دیگر، نگاهی به موارد استفاده از وبگاه ها از سوی گروه های نفرت دارد، مانند شوالیه های

کوکلاکس کلان (1) و اتحاد ملی سازمان نئونازیسم که مقاصد ترغیبی را دنبال می کردند (Duffy 2003). در این مورد، دافی، از روش انتقادی بلاغی به نام تحلیل موضوعی فانتزی استفاده کرد که بر «نظریه همگرایی نمادین» استوار است. همانطور که دافی مشخص کرد، نظریه همگرایی نمادین «نظریه کلی علم معانی است که در آن، گروه ها تصوراتی را درباره گروه و گروه های خارج شکل می دهند و به اشتراک می گذارند، و در نتیجه هویت مشترکی را به وجود می آورند» (صفحه 293) (2) نویسنده، به وسیله تحلیل متون برگزیده موجود در در وبگاه هر گروه توانست گزارش های بلاغی متعددی را شناسایی کند که «ادعای انصاف و عدالت»، «نظم طبیعی و رستاخیز انسان» و «ساکنان اولیه زمین نامیده شدند تا مفهوم جدیدی را به نژاد انتقال دهد» از آن جمله هستند (صفحه 295-305).

4- تحلیل گفتمان

*تحلیل گفتمان (3)

مک کوئیل (2000)، تحلیل گفتمان را استعمال آن «در تمامی اشکال کاربردی زبان و اشکال نوشتاری» می داند، با استناد بر این طرز تفکر که «ارتباط در قالب نوشتار و گفتار رخ می دهد که با موقعیت اجتماعی خاص، موضوع ها، و انواع شرکت کنندگان سازگار می شود» (صفحه 494). تحلیل گفتمانی متون وب، موقعیت اجتماعی - فرهنگی وبگاه ها را در نظر می گیرد؛ یعنی به ساختار معنی محلی آن ها از طریق عناصر گفتاری و دیداری می پردازد. دگربار، متون وب متعلق به گروه های نفرت ظاهر می شوند؛ چنان که بیلینگ (2001) قراردادهای زبانی طنز را همانطور که از سوی وبگاه های «طنز» متعلق به گروه کوکلوکس کلان دیده می شود کشف کرد. او پی برد که در این وبگاه ها از تکذیب کنندگان برای هشدار علیه تفسیر طنزهای نژادپرستانه استفاده می کردند این طنزهای نژاد پرستانه به اقدامات خشونت آمیز فرمان می دادند در حالی که محتوای آن ها برای افراد «شوخی طبع» تعبیری طنزگونه داشت (صفحه 274). هر چند نویسنده به این نتیجه رسید که بخش عمده محتوای اینگونه سایت ها طنزگونه نبوده بلکه از جنس حقیقت است و فراگفتمانی را که او مطرح می سازد «انکار می کند که طنز طنز است» (صفحه 278). متون و تصاویر نژادپرستانه که هم خوانندگان این «طنزها» با گرایش های نژادپرستانه و هم مقاصد اخلاقی اینگونه طنزها را مخاطب قرار می داد شواهدی از گفتمان اهریمنی تری را فراهم می نمود، گفتمانی که بررسی بیلینگ را به تعاملات اجتماعی اوپاش لینچ (4) تشبیه می کند (صفحه 287).

گفتمان در سطح کلان نیز در وب مورد تجزیه و تحلیل قرار گرفته است، و تأکید آن بر وبگاه هایی است که صدای دولت و ملت ها هستند. در بررسی پورسل و کورداس از وبگاه دولت اسلوونی، آن ها با متون موجود در این وبگاه به عنوان واکنشی به خطاهای مشاهده شده در معرفی اسلوونی به عنوان دولت بالکان - منطقه ای که گرفتار جنگ داخلی شد و بنابراین از سوی دولت های غربی به عنوان دولتی

ص: 80

Ku Klux Klan -1

2- فاس (1996)، نظریه همگرایی نمادین را بر اساس دو ادعای اصلی توصیف می کند: «ارتباط واقعیت را می سازد» و «نه تنها واقعیت را برای افراد می سازد بلکه معانی افراد را از نشانه ها می تواند با هم یکی کند تا واقعیت مشترکی را برای شرکت کنندگان به وجود آورند» (صفحه 122).

Discourse analysis -3

Lynch -4

نامطلوب به خاطر اقدامات تروریستی و سرمایه گذاری های تجاری دیگر محکوم شد - برخورد کردند. پورسل و کورداس اظهار می دارند که اسلوانی تلاش می کند تا هویت از دست رفته خود را از طریق متون و تصاویر در وبگاه دولتی احیا کند که نوعی راهکار بلاغی محسوب می شود. اما بر اساس نتیجه گیری آن ها این سایت، سرانجام بخشی از تلاشی مستمر برای مذاکره موقعیت [اسلوانی] در «گفتمان جهانی» است.

شندلر (1998)، بیان کرده است که علاوه بر وبگاه های سازمانی، صفحه های شخصی نیز ممکن است شکلی از گفتمان باشند. شندلر در تحلیل خود از متون وب تولید شده توسط نوجوانان ولز از طریق مصاحبه با سازندگان سایت ها، روش های بسیاری را بررسی کرد که بیان گر تصور نوجوانان از مخاطبان شان بود و درباره مرز بین فضاها ی عمومی و خصوصی زندگی شان صحبت می کردند. بسیاری از نویسندگان جوان وب بیان کردند که سعی دارند تا در سایت های خود با مخاطبانی که احتمالاً علاقه های خود را به اشتراک می گذارند ارتباط برقرار کنند در حالی که افراد دیگر وبگاه را چیزی توصیف می کنند که صرفاً برای خودشان ایجاد کرده اند برخی، دیگر انگیزه ای پیدا می کنند تا بخشی از این اینترنت «پهناور» باشند و عقایدشان را درباره زندگی با دیگران به اشتراک بگذارند.

5- تحلیل دیداری

*تحلیل دیداری (1)

به دلیل قابلیت های چند رسانه ای شبکه جهانی وب، برخی پژوهشگران در صدد برآمده اند تا تحلیل های سنتی از متون وب گفتاری را با تحلیل متون دیداری تکمیل کنند یا بر آن ها غلبه کنند. ایجادکنندگان وبگاه ها نیز از گرافیک ها و تصاویر خلاقانه فراوان با چاشنی سلیقه در وبگاه ها بهره می گیرند، مانند آن چه در JenniCam.org به کار رفته است. از ویژگی های این وبگاه شخصی دوربین 24 ساعته ای با چشم انداز اتاق خواب جنی رینگلی نویسنده وب است. جیمروگلو (1999)، از سایت جنی گم (2) به عنوان بررسی موردی برای کشف مفهوم علم فرمانشی ترکیبی از ماشین و ارگانسیم هاراوای (3) استفاده می کند (صفحه 149). همانند سایبورگ جیمروگلو مشاهده کرد که مرزهای بین بدنه جنی و فناوری از طریق دوربین پیوسته او محو می شود، مثلاً تصویر جنی که در مقابل صفحه نمایش رایانه اش نشسته است. سایت وب گم او تصویر جنی در پشت رایانه اش به نمادی برای آن تبدیل شده است، به طوری که جسمش ذوب و تبدیل به صفحه کلید می شود» (صفحه 441). برای بررسی معنای جنی گم، به عنوان موضوعی جنسیتی و دیداری از نظریه فمینیستی فیلم استفاده شد. تعهد رینگلی در ارائه نگاهی اجمالی به زندگی واقعی او در برگزیده لحظاتی بود که از جلوی وب گم به دلیل برهنگی یا داشتن فعالیت جنسی کنار می رفت. جیمروگلو، مشکل منحصر به فرد این تصاویر را تشخیص داده است: جنی موفق شد تا مرزهای سنتی اختصاص یافته به اندام زن را بین فضاها ی عمومی و خصوصی بشکند، حال آن که با انتقاد برخی فمینیست ها روبه رو شد. آن ها بر این باور هستند که جنی در دام عینیت بخشی اندام زن، که در وب شایع

ص: 81

Visual analysis -1

JenniCam -2

Haraway -3

است، گرفتار شده است. در هر دو مورد «اندام جنی به عنوان کانون معنا، محل کمال، و ریشه معنی واحدی از جنی کم عمل می کند» (صفحه 449).

پژوهشگران دیگر، از تحلیل دیداری برای ارزیابی کاربرد اینترنت از سوی رسانه های خبری برای ارائه و توزیع تصویر استفاده کرده اند. فر پارکز (2001)، تحلیل انتقادی را از تصاویر مربوط به بحران نسل کشی رواندا در سال 1994 ارائه می دهند که رسانه های خبری آمریکا از طریق تلویزیون و اینترنت در دسترس عموم قرار دادند. پوشش زمینی این بحران، که در آن دوربین ها تصاویری را از پناهندگان و حشت زده گرفته بودند و امدادگران غربی سفید پوست به آن ها یاری می رساندند، و نیز تصاویر ماهواره ای از مردمی که دسته جمعی از رواندا می گریختند، با مخالفت روبه رو شدند، زیرا از قربانیان این بحران اختیار را سلب کرده بود. [آن ها] موقعیت سیاسی این واقعه را بیش از اندازه ساده انگاشته بودند و بیننده را از مردم درگیر و مسائل موجود دور نگه می داشتند. فر و پارکز، می افزایند که استفاده از وبگاه ها برای ارائه و پخش تصاویر رواندا بخشی از روند گسترده تر سلطه و مالکیت فناوری های دیداری از سوی رسانه های خبری در غرب است. در واقع، محتوای تولید شده به مصرف مخاطبان غربی یعنی کاربران اصلی اینترنت می رسد. فر و پارکز، به این نتیجه رسیدند که انتخاب تصاویر از سوی رسانه های خبری در پوشش خبری خود از بحران رواندا به همراه پخش گسترده تصاویر از طریق فناوری های رسانه ای غرب دست به دادند تا «آشفتگی» موجود آمریکا را از طریق فرهنگ و سیاست آفریقا تقویت کنند (صفحه 42).

دو مطالعه یادشده، توجه خاصی به تصاویر مورد استفاده در وبگاه ها، به منظور انتقال تفاسیر شان مبذول داشته اند. اما تحلیل های متون دیداری و گفتاری، جدا از منحصر به فرد بودن دو جانبه شان، معمولاً توسط پژوهشگران تلفیقی می گردند. همان گونه که در مطالعه هسکچ و سالیوان، درباره هویت بخشی مدارس از طریق تصاویر و متن ها در وبگاه ها، یا مطالعه و ارنیک درباره وبگاه های مقلد مربوط به مبارزات انتخاباتی و استفاده آن ها از متون گفتاری و غیر گفتاری برای پیشبرد دستور جلسات سیاسی ویژه مشاهده گردید. طراحی دیداری وبگاه ها، از قبیل استفاده از تصاویر و متن ها، و گذاشتن آن ها در صفحه وبگاه، نمونه دیگری از تلفیق این روش هاست (2000 Rivett). بررسی موردی وبگاه «نیوبیتل» ولکس واگن (1)، در مورد صفحه آرایی برگزیده، را که در صفحه اصلی این سایت قرار داشت مورد بحث قرار داد، و نویسنده سایت آن ها را با قراردادهای طراحی مجلات ناپیوسته مقایسه نمود. به علاوه، این سایت، با استفاده از یک پیش زمینه سفید، «هویت دیداری» پیوسته شرکت را منتقل می کند (صفحه 50). تحلیل دیداری وبگاه نیوبیتل مکمل تحلیل متنی است. در صورتی که داستان «حمله بیگانگان» در سرتاسر صفحه هایش به چشم می خورد. سپس، سایت ولکس واگن با طرح دیداری سایت دانشگاهی دانیل شندلر مقایسه گردید. سایت دانشگاهی دانیل شندلر، که مملو از متن بود بر هدف خود به عنوان منبع اطلاعات تأکید می کرد. در این سایت از تصویر کلاسوری به عنوان تصویر زمینه این سایت استفاده شد که بر ماهیت آموزشی آن تأکید می کند.

ص: 82

*قوم نگاری (1)

قوم نگاری، اغلب به عنوان «توصیف و تفسیر گروه یا نظام فرهنگی یا اجتماعی» درک می شود (Creswell 1998P58). تعریف موضوع در پژوهش های قوم نگاری به فرآیند پیچیده ای تبدیل شده است. زیرا در خلال مطالعه فرهنگ های اینترنت و وب انجام می گیرد. جنبه های عملی قوم نگاری، پیش از اینترنت، مطالعه فرهنگ های تک محله ای بوده است که از لحاظ جغرافیایی متمرکز شده اند. با وجود این انسجام پیوسته و پراکندگی فیزیکی ناپیوسته شرکت کنندگان در فرهنگ های اینترنتی مستلزم روشی چند مکانه برای قوم نگاری است. مفهوم جامعه چند مکانه نیز باید در موقعیت های شرکت کنندگان جامعه پیوسته به کار برده شود و ارائه توصیفی تیره از جامعه پیوسته مستلزم روشی متنوع تر از مشاهده ساده گفتمان گروه درون اتاق گفت و گو یا تحلیل کردن محتوای وبگاه های مربوط به آن جامعه است. در عوض، برخی پژوهشگران قوم نگار وب استدلال کرده اند که نیازمند روشی کلی نگر تر هستند، روشی که پژوهشگر تمامی سایت های مربوط به مشارکت و تجربه اعضای جامعه پژوهش را هم به صورت پیوسته و هم ناپیوسته بررسی کند (Howard 2002, Miller 2000, Slater).

میلر و اسلتر (2000) نمونه سودمندی را برای اجرای الگوی چند مکانه از قوم نگاری اینترنتی در مطالعه خود درباره اهالی ترینیداد (2) و اینترنت ارائه کردند. پژوهشگران تالارهای گفت و گو، وب گاه ها، و نیز خانه ها و کافی نت های مردم را بازدید کردند تا از روش های تلافی اینترنت با زندگی های سیاسی، اقتصادی، و مذهبی «ترینی ها» آگاهی یابند. هر سایت، ارائه دهنده چشم انداز منحصر به فردی از آن فرهنگ است. اتاق های گفت و گو مکان هایی هستند که بینشی را از «ترینی» بودن به صورت پیوسته به دست می دهند. وب گاه ها، به عنوان نمایندگان ترینیداد به صورت پیوسته بودند (صفحه 13). میلر و اسلتر، تحلیل محتوایی کیفی را به کار گرفتند تا درباره وب گاه ها به بحث و گفت و گو بپردازند. در حالی که، تلفیق نمادهای ملی نظیر پرچم و تصاویر نقشه ها را که با هویت های شخصی نویسندگان وب سراسر سایت های تحلیل شده اند - متذکر می شوند. افزون بر آن، آن ها به بحث درباره بازتاب دستور جلسات سیاسی در وبگاه ها، با سود آور ترین جاذبه توریستی ترینیداد، یعنی جشنواره کارنیوال، پرداختند که به محور اصلی در این سایت ها تبدیل شده بودند. پژوهشگران علاقه مند به جوامع طرفدار اینترنت، قوم نگاری را روشی سودمند برای مطالعات شان می دانند. بلوستین (2002)، تحلیل های انجام گرفته از وب گاه های طرفدار و انجمن های مربوط را در ترسیم خود از تعصب، و نمایش تلویزیونی «بافی»، قاتل خون آشام» گنجانند. او با استفاده از تحلیل های خود از متون وب و اطلاعات گردآوری شده از مصاحبه های شخصی و بازدید از اتاق های طرفداران توانست به نتیجه گیری هایی برسد که چگونه طرفداران بافی (که اغلب نوجوان بودند) بین جلوه فانتزی و جلوه های سحرآمیز نمایش تمایز قائل می شوند تا تخیل و «تفریح» جوانی را حفظ کنند، مادامی که به فضای جدی و ترسناک ابهام اخلاقی موجود در بزرگسالی وارد می شوند (صفحه 440).

ص: 83

Ethnography -1

2- ترینیداد بزرگ ترین جزیره کشور ترینیداد و توباگوست. این جزیره جنوبی ترین جزیره دریای کارائیب است که تنها 11 کیلومتر با کرانه های ونزویلا فاصله دارد. ترینیداد 4,768 کیلومتر مربع مساحت دارد.

*تحلیل شبکه (1)

به عقیده گارتون و همکارانش (1997)، «شبکه اجتماعی، به مجموعه افراد (یا سازمان ها یا نهادهای اجتماعی دیگر) اطلاق می شود که از طریق مجموعه ای از روابط اجتماعی نظیر دوستی، همکاری یا تبادل اطلاعات به هم مرتبط می گردند» (پاراگراف 2). هنگامی که این روابط اجتماعی با استفاده یا ایجاد وبگاه ها بنا برقرار یا ایجاد می گردند، فرایندها وسیله ای را برای مطالعه الگوهای ارتباطی میان اعضای شبکه ارائه می دهند. پارک (2003)، این نوع تحلیل در پژوهش وب را تحلیل فرایوند می داند که وبگاه ها را نمایان گر افراد، گروه ها، سازمان ها، و دولت ها و ملت ها، و فرایندهای بین سایت ها را «پیوند[های] رابطه ای» تصور می کنند (صص 50-51). برای مثال، هالویس (2000)، به منظور درک بهتر این موضوع که آیا وب به درستی شبکه «جهان گستر» سایت ها را رواج می دهد یا اینکه مرزهای ملی در حال پیوستن به مرزهای اینترنت هستند؟، به مطالعه فرایندهای موجود در وب گاه ها پرداخت. او با مطالعه 4000 وبگاه به بررسی پیوندهای آن وبگاه ها با صفحه های خارج از آن وبگاه ها پرداخت که از هویت میزبان ملی وبگاه هایی که به آن ها پیوند می یابند، پیروی می نمایند. هالویس، با استفاده از این فرآیند اینگونه نتیجه گرفت که در حقیقت بیشتر وبگاه های موجود مورد بررسی اصولاً به سایت های درون فرهنگ های ملی مرتبط می شوند. با وجود این، در مقایسه با انواع دیگر فناوری های اطلاعاتی مانند پست زمینی و تلویزیون او بر این باور بود که اینترنت بین المللی ترین فناوری اطلاعاتی است، و نسبت به دیگر رسانه ها در داخل و خارج از آمریکا در سراسر مرزهای بین المللی از تعداد مراجعات بیشتری برخوردار است (صفحه 23).

همان طور که پیشتر ذکر شد، تحلیل شبکه در سطح سازمانی، همزمان با تحلیل محتوایی در مطالعه وبگاه های ضدجهانی سازی، برای درک بهتر این مسئله مورد استفاده قرار گرفت که این سایت ها تا چه حد با هم تلفیق شده اند (Van Alst and Walgrave 2002). پیوندهای موجود بین 17 سایت، با استفاده از نرم افزار دیداری سازی شبکه (Pajek: <http://vlado.fmf.uni-lj.si/pub/networks/pajek>)، مورد تجزیه و تحلیل قرار گرفت. نرم افزار یاد شده، نقشه گرافیکی پیوندهای میان وبگاه ها را ایجاد می کرد. سایت های بسیار مرتبط و دارای ارجاعات بسیار با گردآوری پیوندها در موقعیت خودشان، و خروج از آن ها به آسانی شناسایی شدند. حال آن که، سایت های پراکنده تنها به وسیله یک یا دو پیوند با دیگر سایت های شبکه پیوند می یابند، دور افتاده به نظر می رسند. ون آلست و والگریو دریافتند که سایت های ضد جهانی سازی تا اندازه ای یکپارچه شده اند، اما یادآور می شوند که ارزیابی ماهیت روابط بین سازمان هایی که صرفاً مبتنی بر موجودیت پیوندهای بین سایت ها می باشند، دشوار است.

در موارد دیگر، فرایندها به منظور کشف شبکه های اجتماعی و اطلاعاتی افراد پیوسته مورد تجزیه و تحلیل قرار گرفته اند بیشتر مرورگرهای وب، مانند اینترنت اکسپلورر (3) و Netscape Navigator، مجهز

به کار کرد «تاریخچه (1)» هستند که نشانی های اینترنتی را که شخص در هنگام بازدید از سایت های شبکه به آن ها دسترسی داشته است، ذخیره می کنند. تاچر و گرینبرگ (1997)، از این کار کرد برای تحلیل رفتارهای مرورگرانه 23 نفر طی یک دوره 6 هفته ای استفاده کردند. آن ها به دنبال الگوهایی در بازدید صفحه های وب بودند. این دو پژوهشگر، متوجه شدند که تقریباً دو سوم بازدیدهای صفحه های وب از سوی آزمودنی ها، صفحه هایی هستند که قبلاً بازدید شده بودند (صفحه 112). همچنین، تعداد صفحه هایی که اغلب توسط مخاطبان دوباره بازدید می شوند، نسبتاً کم بودند. تاچر و گرینبرگ، فرض می کنند که کارکرد «بازگشت (2)» در مرورگرهای وب، که 30 درصد از فعالیت های مرورگرانه آزمودنی ها را تشکیل می دهد، ممکن است عامل مؤثری برای علاقه کاربران اینترنت برای بازدید دوباره صفحه ای واحد در وب باشد (صفحه 131).

علاوه بر شبکه های اجتماعی مطالعه شبکه های مدارک، به دلیل بزرگی شبکه وب رو به فزونی رفته است، و این بدان معناست که اشکال گوناگون تحلیل علمی از شبکه اجتماعی امکان پذیر است. تحلیل شبکه های مدارک می تواند فرصت هایی را برای فراتحلیل محتوا به وجود آورد. هنزینگر و لورنس (2004)، درباره روش های نمونه برداری از صفحه های وب به بحث و گفت و گو می پردازند تا «به طور خودکار نمونه وسیعی از علایق و فعالیت های انجام گرفته در این دنیا را ... به وسیله تجزیه و تحلیل ساختار پیوند وب و چگونگی روی هم انباشته شدن پیوندها در طول زمان تجزیه و تحلیل کنند» (صفحه 5186). آیزن و همکارانش (2004)، بر این ادعا هستند که استفاده از داده ها در وبگاهی با ترافیک بالا می تواند اطلاعاتی را درباره وقایع خارجی و احساسات ناگهانی در دسترس عموم قرار دهد که ممکن نیست به تنهایی از تحلیل محتوا و ساختار پیوند دسترس پذیر شود (صفحه 5254).

8- ملاحظات اخلاقی

اهمیت نیاز به پژوهش با رعایت اصول اخلاقی برای پژوهشگران اینترنت پوشیده نیست. اما خط مشی ها و ابزارهای پژوهش اخلاقی به روشنی مشخص نیستند. از دیدگاه پژوهشگران وب، و به ویژه آن هایی که مواد آرشیوی وب را استخراج می کنند، ملاحظات اخلاقی به صورت ضرورت در جریان پژوهش نمود می یابند. برخی اطلاعات موجود در وب و اطلاعات در دست آرشیو، که محرمانه به شمار می آیند، ممکن است سهواً دسترس پذیر شده، و سپس توسط موتورهای کاوشی نظیر گوگل نمایه و ذخیره شوند. دیدگاه های مربوط به امکان استفاده از اطلاعات خصوصی، که به صورت عمومی آرشیو شده اند، متعدد و خارج از حوصله این مقاله می باشند. با وجود این، باید خاطرنشان کرد که تمامیت موتورهای کاوش در جست و جوی وب می تواند منجر به ایجاد آرشیوهای اطلاعاتی شود که کاربران عموماً نه قصد فاش کردن آن ها را دارند و نه آن ها را آرشیو کرده اند. در حالی خوش بینانه، برآوردی سرعتی که موتورهای کاوش وب را نمایه سازی می نمایند، نشان داد که سرعت نمایه سازی و آرشیوسازی گوگل از ایجاد صفحه های جدید

ص: 85

وب فراتر رفته است (Whelan 2004). احتمالاً مسائلی که پژوهشگران حوزه مطالعات اینترنت مدت ها به بحث و تبادل نظر گذاشته اند، و نیز مطالعات انجام شده بر روی ارتباطات با واسطه رایانه، با توجه به تقابل داده های عمومی با داده های خصوصی، در آینده بسیار نزدیک در میان پژوهشگران حوزه مطالعات وب و حرفه مندان حوزه آرشیو وب و دیگر مواد الکترونیکی از اهمیت بالایی برخوردار خواهند شد.

9- نتیجه گیری

پژوهش های حوزه وب به وضوح سودمندی خود را نشان داده اند. این پژوهش ها بدون آن که در سایه تحلیل های متنی جوامع مبتنی بر متن قرار داشته باشند، ثابت کرده اند که مفاهیمی نظیر جامعه، فرهنگ، رفتار و ساختارهای معنایی می توانند در اینترنت به طور مؤثری بررسی شوند. در حالی که غنا و گوناگونی پژوهش های وب نوید بخش است، هنوز مسائل زیادی وجود دارد که باید کشف شوند.

میراث تحلیل متنی به جامانده از ارتباطات با واسطه رایانه در سوگیری شدید نسبت به تحلیل های زبانی در مطالعات حوزه وب مشهود است. «مطالعات مفقود» به مطالعاتی اطلاق می شود که وب را مانند چشم اندازی چندرسانه ای بررسی می کند، و توصیف می کند که چگونه تصاویر و صدا برای از میان برداشتن ارتباطات گفتاری مورد استفاده قرار می گیرند تا بر موانع زبانی در رسانه جهانی فائق آیند. به استثنای اشاره به کاربرد صدا در وبگاه های ترینیدادی، در مطالعه قوم نگاری میلر و اسلتر (2000)، وب گاه های شنیداری کاملاً از بخش کنونی پژوهش وب غایب بودند. احتمال دارد که پروژه های آینده به کاربردهای بلاغی صدا در وب از طریق تحلیل های وب گاه های تجاری یا سیاسی پردازند. دیگر مطالعاتی که به حوزه چند رسانه ای ها گرایش دارند توانستند روش هایی را که از صدا و تصاویر برای ایجاد هویت افراد موجود در اینترنت از افراد حقیقی و سازمان ها گرفته تا دولت های ملی استفاده می شدند، بررسی کنند.

به منظور افزایش تعداد پژوهشگران این حوزه دسترسی به همه انواع محتوای وب مربوط به تمامی دوره های زمانی، به ویژه بازآفرینی پیوندهای درون اشیا و محتوایی و میان وبگاه ها الزامی است. به علاوه، باید مرورگرها و ابزارهای دیگری در دسترس قرار گیرند که محتوای وب به وسیله آن ها رصد شود تا بتوان تجربه کاربر را به بهترین شکل درک کرد. بنابراین، چالش ها فراتر از حفظ و نگهداری محتوا رفته، و شامل حفظ ساختار و رویارویی با محتوا نیز می شود.

منابع

Aizen, J., Huttenlocher, D., Kleinberg, J., Novak, A. (2004). Traffic-based feedback on the Web. . 1
Proceedings of the National Academy of Sciences of the United States of America, 101(1), 5254-5260

Billig, M. (2001). Humour and hatred: The racist jokes of the Ku Klux Klan. Discourse .2

- Bloustein, G. (2002). Fans with a lot at stake: Serious play and mimetic excess in Buffy the Vampire Slayer. *European Journal of Cultural Studies*, 5(4), 427–449 .3
- Carswell, A. D. Venkatesh, V. (2002). Learner outcomes in an asynchronous distance education environment. *International Journal of Human–Computer Studies*, 56, 475–494 .4
- Chandler, D. Roberts–Young, D. (1998). The construction of Identity in the Personal Homepages of Adolescents. Retrieved April 15, 2002 from <http://www.aber.ac.uk/media/Documents/short/strasbourg.html> .5
- Creswell, J. W. (1998). *Qualitative inquiry and research design*. Thousand Oaks, CA: Sage .6
- Duffy, M. E. (2003, July). Web of hate: A fantasy theme analysis of the rhetorical vision of hate groups online. *Journal of Communication Inquiry*, 27(3), 291–312 .7
- Ess, C. (2002). Ethical decision making and Internet research: Recommendations from the AoIR working committee. Association of Internet Researchers. Retrieved July 12, 2004, from <http://www.aoir.org/reports/ethics.pdf> .8
- Fair, J. E. Parks, L. (2001). Africa on Camera: Television news coverage and aerial imaging of Rwandan refugees. *Africa Today*, 48(2), 34–57 .9
- Foss, S. K. (1996). *Rhetorical criticism: Exploration and practice*. Prospect Heights, IL: Waveland .10
- Garton, L., Haythornthwaite, C., Wellman, B. (1997). Studying online social networks. *Journal of Computer–Mediated Communication*, 3 (1). Retrieved April 4, 2004 from <http://www.ascusc.org/jcmc/vol3/issue1/garton.htm> .11
- Geertz, C. (1973). *The interpretation of cultures*. New York: Basic Books Graphic, Visualization and Usability (GVU) Center (1998). 10th WWW User Survey. Atlanta, GA: Georgia Institute of Technology. Retrieved March 25, 2004 from <http://www.cc.gatech.edu/gvu/user-surveys/survey-1998-10> .12
- Halavais, A. (2000). National Borders on the World Wide Web. *New Media Society*, 1(3), 7–28 .13
- Haraway, D. (1991). *Simians, cyborgs, and women: The reinvention of nature*. New York: Routledge .14

Henzinger, M. Lawrence, S. (2004). Extracting knowledge from the World Wide Web. Proceedings of . 15
,the National Academy of Sciences of the United States of America

ص: 87

- Hesketh, A. J. Selwyn, N. (1999). Surfing to school: The electronic reconstruction of institutional . 16
identity. *Oxford Review of Education*, 25(4), 501-520
- Howard, P. N. (2002). Network Ethnography and the Hypermedia Organization: New Media, New . 17
Organizations, New Methods. *New Media and Society*, 4(4), 550-574
- Jimroglou, K. M. (1999). A Camera with a view: JenniCam, visual representation, and . 18
cyborgsubjectivity. *Information, Communication and Society*. 2(4), 439-453
- Kroeber, A. (2001). Postmodernism, Resistance, and Cyberspace: Making Rhetorical Spaces for . 19
Feminist Mothers on the Web. *Women's Studies in Communication*, 24(2), 218-240
- Lee, W., Hong, J., Lee, S. (2003). Communicating with American consumers in the post 9/11 climate: .20
An empirical investigation of consumer ethnocentrism in the United States. *International Journal of
Advertising*, 22, 487-510
- McQuail, D. (2000). *McQuail's mass communication theory* (4th ed.). London, UK: Sage .21
- Miller, D. Slater, D. (2000). *The Internet: An ethnographic approach*. Oxford, UK: Berg .22
- Milne, G. R. Culnan, M. J. (2002). Using the content of online privacy notices to inform public policy: .23
A longitudinal analysis of 1998-2001 US Web surveys. *The Information Society*, 18, 345-359
- O'Neill, E. T., Lavoie, B. F., Bennett, R. (2003). Trends in the evolution of the public Web, 1998-2002. .24
D-Lib Magazine, 9(4). Retrieved March 29, 2004 from <http://wcp.oclc.org>
- Park, H. W. (2003). Hyperlink network analysis: A new method for the study of social structure on the .25
Web. *Convergence*, 25(1), 49-61
- Pitts, M. J. Harms, R. (2003). Radio websites as promotional tools. *Journal of Radio Studies*, 10(2), .26
270-282
- Purcell, D. Kodras, J. E. (2001). Information technologies and representational spaces at the outposts of .27
the global political economy. *Information, ommunication and Society*, 4(3), 341-369
- Rivett, M. (2000). Approaches to analyzing the Web text: A consideration of the Web site as an .28

emergent cultural form. *Convergence*, 6(3), 34–56

?Sarraf, S. (1999). A survey of museums on the Web: Who uses museum Websites .29

ص: 88

- Silver, R. C., Holman, E. A., McIntosh, D. N., Poulin, M., Gil-Rivas, V. (2002). Nationwide . 30 longitudinal study of psychological responses to September 11. *Journal of the American Medical Association*, 288(10), 1235-1244
- Spooner, T. (2002). Internet use by region in the United States. Pew Internet American Life Project. . 31 Washington, DC. Retrieved March 25, 2004 from <http://www.pewinternet.org/reports/pdfs/PIP-Regional-Report-Aug-2003.pdf>
- Tauscher, L. Greenberg, S. (1997). How people revisit Web pages: Empirical findings and implications . 32 for the design of history systems. *International Journal of Human- Computer Studies*, 47, 97-137
- Van Aelst, P. Walgrave, S. (2002, December). New media, New movements? The role of the Internet in . 33 shaping the 'anti-globalization' movement. *Information Communication and Society*, 5(4), 465-493
- Warnick, B. (1998). Appearance or reality? Political parody on the Web in Campaign '96. *Critical . 34 Studies in Mass Communication*, 15, 306-324
- Whelan, D. (2004, 16). Google Me Not. *Forbes*, 174(3), 102-104 . 35

اطلاع از مشخصات وب ملی یکی از نیازمندی‌های اصلی کشور در رابطه با سیاست گذاری در حوزه فناوری اطلاعات است. از آن جا که متأسفانه تا به حال، تحقیقی در این خصوص، انجام نشده است، نیاز مبرمی برای پژوهش و کسب اطلاع از وضعیت فعلی وب فارسی، احساس می‌شود. با درک این ضرورت، در این مقاله نتایج حاصل از مطالعات صورت گرفته در مورد خصوصیات و ویژگی‌های وب فارسی مطرح خواهد شد. برای این کار سامانه نیمه خودکاری طراحی و پیاده سازی شده است. هدف از این سیستم، استخراج شاخص‌های مختلف از قبیل حجم محتوای فارسی، تنوع محتوا و نیز عمر محتوا می‌باشد. از این آمار می‌توان جهت هدفمند نمودن برنامه‌های آتی در خصوص ساماندهی وب فارسی و نیز ارائه راهکارهایی برای حل معضلات فعلی، بهره گرفت. آنالیزهای انجام شده توسط سامانه فوق‌الذکر بر روی حدود یازده هزار وبگاه ثبت شده در دامنه IR، مشتمل بر اغلب سازمان‌ها و ارگان‌های دولتی، وزارتخانه‌ها، شرکت‌ها و دانشگاه‌ها است که بالغ بر حدود دو میلیون صفحه می‌باشند.

کلید واژه: دامنه IR، وب فارسی، مشخصات وب

*خصوصیات وب ایران: مریم پیروزمند (1)

1. مقدمه

توسعه و رشد نمایی وب باعث شده است تا با حجم عظیمی از اطلاعات شامل اسناد با فرمت های متفاوت از جمله متن، صوت، تصویر و غیره در مکان ها و سازمان های مختلف مواجه شویم. در عین حال، گسترش روز افزون نیازمندی به وب که زاینده فزونی عرضه خدمات از طریق وب است، باعث شده تا همواره کاربران بیشتری به استفاده از وب، راغب شوند. بنا بر آمارهای موجود در حال حاضر بالغ بر 130 میلیارد صفحه در محیط وب از طریق جویشرهای مطرح دنیا قابل دسترسی و جستجو هستند [1]. به دلایل مختلف از جمله شروع توسعه وب توسط کاربران انگلیسی زبان، غالب سرویس ها و خدمات عرضه شده از طریق وب بویژه سرویس های جستجو، بصورت انگلیسی ارائه می شوند. گرچه در طی سال های اخیر، موتورهای جستجوی عمده ای مانند Yahoo و Google، سرویس جستجو را به زبان های دیگر نیز عرضه کرده اند، اما متأسفانه به دلایل مختلف بالاخص تحریم های سیاسی و اقتصادی، این خدمات هیچگاه به حوزه زبان فارسی گسترش پیدا نکرده است. از سوی دیگر، در داخل کشور نیز به علت عدم وجود شناخت کافی از مختصات وب فارسی موتورهای جستجوی توانمندی بوجود نیامده اند.

بطور کلی، شناخت مختصات وب ملی یک کشور، نمایان گر شاخص های توسعه یافتگی در آن کشور

ص: 91

نیز محسوب می شود. به عنوان مثال، در کشورهای توسعه یافته، حجم وب ملی شامل تعداد وبگاه های مختلف و همچنین تنوع خدمات قابل عرضه در این رسانه، بسیار زیاد است. همچنین کیفیت و به روز بودن اطلاعات عرضه شده نیز در دامنه وب این قبیل کشورها در مقایسه با دیگر کشورها بسیار بهتر است.

از سوی دیگر، با توجه به اینکه زیر ساخت اصلی اشاعه دولت الکترونیک، محیط وب می باشد، لذا مطالعه و شناسایی ویژگی های این محیط، کمک شایانی به انجام برنامه ریزی مناسب جهت تدوین سیاست های اجرایی جهت تحقق اهداف دولت الکترونیک طی مراحل مختلف، خواهد کرد. با درک این مطلب، در این مقاله سعی می شود تا به ویژگی های اصلی وب ایران پرداخته شود. شاخص های مورد مطالعه، به گونه ای انتخاب شده است تا بتوان از آن ها در جهت بهبود عملکرد سیستم های جستجوی وب، بهره گرفت.

با توجه به مسائل فوق در این مقاله سعی می شود تا حد ممکن، محتوای وب فارسی کشور را به صورت خودکار، تحلیل و ارزیابی کرد (ابزار نظارت خودکار) تا در مقاطع زمانی مختلف، بتوان عمل ارزیابی را با کمترین هزینه، تکرار نمود. لازم به ذکر است که در این مقاله حدود 11 هزار وبگاه با پسوند IR، که شامل تقریباً دو میلیون صفحه است پردازش و تحلیل شده اند. البته تعداد وبگاه های فارسی بیش از این مقدار است اما به دلیل دشواری یافتن وبگاه های فارسی زبان، در این پژوهش به وبگاه های ثبت شده در دامنه IR بسنده شده است. علت این دشواری این است که همانطور که در قسمت نتایج، به تفصیل ذکر خواهد شد. متأسفانه دامنه IR به عنوان تنها دامنه فارسی زبان نیست و بسیاری از وبگاه های فارسی زبان، در دامنه های مختلف و حتی غیر مرتبطی از قبیل .com, net, org و دامنه های دیگر ثبت شده اند.

همانگونه که پیش از این ذکر شد، هدف اصلی این مطالعه، تعیین وضعیت فعلی وب فارسی است.

بر این اساس، هدف نهایی، استخراج شاخص های زیر از "وب ایران" می باشد:

1. پاسخگویی به سؤالات کلی درباره وب ایران از جمله:

1-1. درصد پیوندهای معتبر چقدر است؟ تعیین این شاخص بطور ضمنی، نرخ تغییرات وب را نیز تعیین خواهد کرد و زمان بندی برای خزشگر وب را دقیق خواهد کرد.

2-1. توزیع صفحات وب فارسی از لحاظ محتوا، چگونه است؟ یعنی چه حجمی از محتوای وب در کلاس های علمی، تجاری، روزنامه، خبر، وبلاگ و غیره، قابل طبقه بندی است؟

3-1. ساختار محتوایی صفحات وبگاه ها چگونه است؟ بطور مشخص می خواهیم بدانیم که:

• درصد کدینگ های زبانی مختلف استفاده شده نظیر Windows- 1252, Windows - 1256 و UTF-8 به چه صورت است؟ بدین ترتیب مشخص خواهد شد که تبدیل کدینگ ها در سیستم های بازایی وب چقدر حائز اهمیت است؟

• چه کسری از صفحات، عنوان (1) مناسب دارند؟ وجود صفحات با عنوان مناسب و گویا، موجب بهبود کیفیت بازایی توسط سیستم و نیز سهولت دسترسی توسط کاربران خواهد شد.

1-4. فایل های غیر متنی مانند Doc PPT PDF و Image چند درصد از صفحات را تشکیل می دهند؟

1-5. نرخ به روزآوری، تغییر، ایجاد (عمر صفحه) چقدر می باشد؟ استخراج این اطلاعات، نرخ تغییرات وب را تعیین خواهد کرد و موجب برنامه ریزی و زمان بندی برای خزش گر وب را فراهم خواهد آورد.

1-6. سرعت دسترسی به وبگاه چقدر است؟ این آگاهی نیز در تنظیم عملکرد خزشگر وب موثر خواهد بود.

1-7. تعداد دسترسی به وبگاه ها و صفحات در مقاطع زمانی مختلف یا عبارت دیگر الگوی دسترسی به وبگاه ها و صفحات چگونه می باشد؟

1-8. تعداد کل صفحات فارسی، میانگین تعداد صفحات هر وبگاه و حجم آن ها چقدر است؟

1-9. تعداد لغاتی که در تمام صفحات محاسبه شده چقدر می باشد؟

1-10. چند درصد صفحات شامل هر دو محتوای فارسی و انگلیسی است؟

لازم به ذکر است در این فاز به دلایلی مانند پیچیدگی کار استخراج شاخص های مربوط به فازهای تراکنش و تبدیل به کارهای آینده موکول شده است. نتایج این مقاله را می توان در تعریف و پیاده سازی پروژه هایی مانند موتور جستجوی ملی و درگاه دولتی استفاده کرد. به علاوه این پروژه در تدوین راهکارهای آینده جهت تحقق سریع دولت الکترونیک در کشور مفید فایده خواهد بود.

2. کارهای مرتبط در داخل و خارج

تاکنون در داخل کشور مکانیزم ارزیابی وب بدین صورت انجام نشده است، اما بعضی از کشورها کاری شبیه به این پروژه را انجام داده اند. برای مثال در تایلند پروژه ای تحت عنوان "ابزار نظارت خودکار بر پروژه دولت الکترونیک تایلند" [2] انجام شده است. هدف این کار استخراج تمام شاخص های دولت الکترونیک از وب تایلند می باشد. بر اساس نتایج این مطالعه که حاصل بررسی حدود 150 وبگاه دولت تایلند در سال 2002 است، حدود 31% این سایت ها، صرفاً به ارایه اطلاعات می پردازند، حدود 57% امکان تعامل محدود کاربران را با وبگاه فراهم می کنند و تنها حدود 11%، بستر لازم برای اجرای تراکنش های مورد نیاز کاربران را در اختیار وی قرار می دهند.

به صورت مشابه، کارهایی برای استخراج مشخصه های وب در کشورهای اسپانیا [3]، کره جنوبی [4]، استرالیا [5]، پرتغال [6]، اروپا [7]، [8] و آرژانتین [9] انجام شده است. در کنار این تحقیقاتی نیز در مورد مطالعه ویژگی های کلی وب صورت گرفته است [10]، [11]، [12] در فعالیت های فوق بیشتر وب کشورها از دیدگاه ساختاری و شکل گراف وب مورد بررسی قرار گرفته است و به علاوه پارامترهایی مانند توزیع اندازه صفحات و وبگاه ها، نرخ بروزآوری آن ها، رتبه آن ها در موتورهای جستجو و عمر صفحات مورد بررسی قرار گرفته ولی از منظر محتوا و سرویس های ارائه شده کاری انجام نگرفته است. در این مقاله علاوه بر استخراج پارامترهای فوق برای وب فارسی، از دید دولت الکترونیک نیز به وب

ایران توجه گردیده است.

در تحقیق صورت گرفته در مرجع [10]، ساختار کلی گراف وب جهانی مورد مطالعه قرار گرفته است. بررسی های صورت گرفته حاکی از وجود ساختاری موسوم به مدل پایونی (1) در مورد گراف وب است. بعنوان مثال، طی بررسی انجام شده در سال 2000 میلادی روی حدود دویست میلیون صفحه وب، همانگونه که شکل شماره یک نشان می دهد، گراف وب، شامل چهار بخش کلی است که عبارتند از:

الف- بخش هسته مرکزی (2) که حدود 25% کل وب را تشکیل می دهد.

ب- بخش ورودی (3) که لینک هایی به بخش هسته مرکزی دارد و حدود 25% کل وب را تشکیل می دهد.

ج- بخش خروجی (4) که لینک هایی از بخش هسته مرکزی دارد و حدود 25% کل وب را تشکیل می دهد.

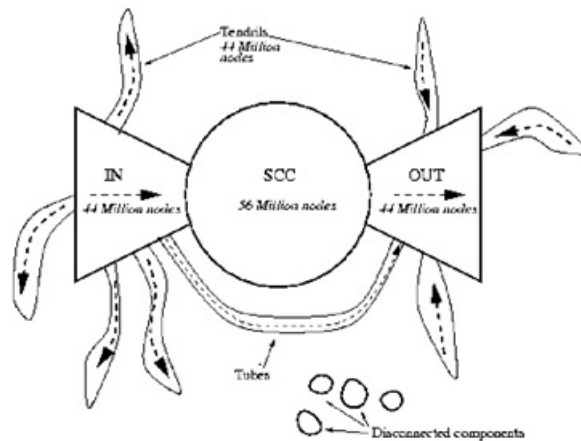
د- بخش های پراکنده که یا به صورت جزایر پراکنده (و بعضاً به سختی قابل دسترس) هستند یا واسطه هایی بین بخش های ورودی و خروجی هستند. این قسمت نیز حدود 25% کل وب را تشکیل می دهد.

عکس

ایران توجه گردیده است.

در تحقیق صورت گرفته در مرجع [۱۰]، ساختار کلی گراف وب جهانی مورد مطالعه قرار گرفته است. بررسی‌های صورت گرفته حاکی از وجود ساختاری موسوم به مدل پایونی^۱ در مورد گراف وب است. بعنوان مثال، طی بررسی انجام شده در سال ۲۰۰۰ میلادی روی حدود دویست میلیون صفحه وب، همانگونه که شکل شماره یک نشان می‌دهد، گراف وب، شامل چهار بخش کلی است که عبارتند از: الف-بخش هسته مرکزی^۲ که حدود ۲۵٪ کل وب را تشکیل می‌دهد.

ب-بخش ورودی^۳ که لینک‌هایی به بخش هسته مرکزی دارد و حدود ۲۵٪ کل وب را تشکیل می‌دهد. ج-بخش خروجی^۴ که لینک‌هایی از بخش هسته مرکزی دارد و حدود ۲۵٪ کل وب را تشکیل می‌دهد. د-بخش‌های پراکنده که یا به صورت جزایر پراکنده (و بعضاً به سختی قابل دسترس) هستند یا واسطه‌هایی بین بخش‌های ورودی و خروجی هستند. این قسمت نیز حدود ۲۵٪ کل وب را تشکیل می‌دهد.



شکل ۱. مدل پایونی گراف وب

۳. سامانه خودکار ارزیابی وب ایران

در این قسمت، ابتدا معماری سیستم ارزیاب خودکار، تشریح می‌شود. این سامانه، نسخه گسترش یافته ابزار استفاده شده در مرجع [۱۳] است که با هدف بررسی مشخصات و بگانه‌های ثبت شده در دامنه IR اعم از دولتی و غیر دولتی می‌باشد. شکل شماره دو شمای کلی این سامانه و تعامل اجزای آن را با

-
1. Bow-Tie model
 2. Central core
 3. In
 4. Out

شکل ۱. مدل پایونی گراف وب

۳. سامانه خودکار ارزیابی وب ایران

در این قسمت، ابتدا معماری سیستم ارزیاب خودکار، تشریح می‌شود. این سامانه، نسخه گسترش یافته ابزار استفاده شده در مرجع [13] است که با هدف بررسی مشخصات و بگانه‌های ثبت شده در دامنه IR اعم از دولتی و غیر دولتی می‌باشد. شکل شماره دو شمای کلی این سامانه و تعامل اجزای آن را با

Bow-Tie model -1

Central core -2

In -3

Out -4

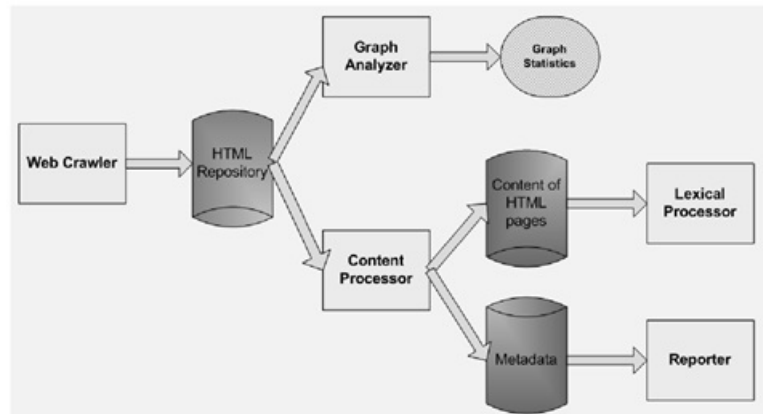
یکدیگر، نشان می دهد.

همان گونه که در شکل شماره دو دیده می شود، ابتدا یک خزشگر وب، بر اساس یک لیست اولیه از وبگاه های ثبت شده در دامنه IR که از قبل تهیه شده است و با توجه به پارامترهای تنظیم عملکرد آن از قبیل عمق خزش، حداکثر صفحات وبگاه و غیره، خزش را با توجه به گراف حاصل از پیوندها، انجام می دهد و صفحات بدست آمده را در یک مخزن موقت، شاخص بندی و ذخیره می کند. در مرحله بعد، این مخزن توسط واحد تحلیل گر گراف وب، مورد بررسی قرار می گیرد و آمارهای مختلفی را مورد ویژگی های گراف متناظر از قبیل قطر گراف، متوسط فاصله بین هر دو گره و غیره، استخراج می کند.

عکس

یکدیگر، نشان می‌دهد.

همانگونه که در شکل شماره دو دیده می‌شود، ابتدا یک خزشگر وب، بر اساس یک لیست اولیه از وبگاه‌های ثبت شده در دامنه IR که از قبل تهیه شده است و با توجه به پارامترهای تنظیم عملکرد آن از قبیل عمق خزش، حداکثر صفحات وبگاه و غیره، خزش را با توجه به گراف حاصل از پیوندها، انجام می‌دهد و صفحات بدست آمده را در یک مخزن موقت، شاخص‌بندی و ذخیره می‌کند. در مرحله بعد، این مخزن توسط واحد تحلیل گر گراف وب، مورد بررسی قرار می‌گیرد و آمارهای مختلفی را مورد ویژگی‌های گراف متناظر از قبیل قطر گراف، متوسط فاصله بین هر دو گره و غیره، استخراج می‌کند.



شکل ۲. معماری کلی سامانه خودکار ارزیابی وب ایران

از سوی دیگر، واحد تحلیل گر محتوا با استفاده از تجزیه کننده HTML، صفحات این مخزن را بررسی نموده و محتوا و داده‌های توصیفی آنها را استخراج می‌کند. این قسمت که توسط زبان جاوا پیاده‌سازی شده است، بصورت مستقل و از طریق پارامترهای قابل تنظیم، اجرا می‌شود. ورودی این بخش، کل صفحات HTML است که شامل برچسب‌های مختلف و نیز نویسه‌های موجود در صفحه می‌باشد. شایان ذکر است که در هنگام جمع‌آوری صفحات از وب، توسط خزشگر وب بصورت خودکار به هر صفحه، یک شناسه عددی نسبت داده می‌شود که آن صفحه را بصورت یکتا مشخص می‌کند. جدول شماره یک، نمونه‌ای از ورودی‌های تجزیه کننده را نشان می‌دهد.

محتوای استخراج شده از صفحات، برای پردازش در اختیار واحد تحلیل گر واژگان، قرار داده می‌شود تا کلمات فارسی را از آن استخراج کند و در یک Lexicon ذخیره نماید. همچنین از اطلاعات بدست آمده می‌توان برای استخراج آمار متنوع از کلمات استفاده کرد.

1. Parser
2. Script

شکل ۲. معماری کلی سامانه خودکار ارزیابی وی ایران

از سوی دیگر، واحد تحلیل گر محتوا با استفاده از تجزیه کننده (1) Parser صفحات این مخزن را بررسی نموده و محتوا و داده‌های توصیفی آن‌ها را استخراج می‌کند. این قسمت که توسط زبان جاوا پیاده‌سازی شده است، بصورت مستقل و از طریق پارامترهای قابل تنظیم اجرا می‌شود و ورودی این بخش کل صفحات HTML است که شامل برچسب‌های مختلف و نیز نویسه (2) های موجود در صفحه می‌باشد. شایان ذکر است که در هنگام جمع‌آوری صفحات از وب توسط خزشگر وب بصورت خودکار به هر صفحه، یک شناسه عددی نسبت داده می‌شود که آن صفحه را بصورت یکتا مشخص می‌کند. جدول شماره یک، نمونه‌ای از ورودی‌های تجزیه کننده را

نشان می دهد.

محتوای استخراج شده از صفحات برای پردازش در اختیار واحد تحلیل گر واژگان، قرار داده می شود تا کلمات فارسی را از آن استخراج کند و در یک Lexicon ذخیره نماید. همچنین از اطلاعات بدست آمده می توان برای استخراج آمار متنوع از کلمات استفاده کرد.

ص: 95

HTML -1

Script -2

بخش مهم دیگر این معماری، واحد تهیه گزارش است. داده های توصیفی استخراج شده از صفحات ورودی به عنوان ورودی به این واحد، ارسال می شود تا مورد بررسی قرار گیرد. در نتیجه این فرآیند آمارهای ارزشمند مختلفی نظیر درصد استفاده از کدینگ های مختلف توزیع صفحات و وبگاه ها در طبقه بندی های مختلف سرویس های ارائه شده نظیر جستجو، امنیت و غیره بدست می آید.

عکس

۹۶ مدیریت منابع اطلاعاتی وب

بخش مهم دیگر این معماری، واحد تهیه گزارش است. داده های توصیفی استخراج شده از صفحات ورودی به عنوان ورودی به این واحد، ارسال می شود تا مورد بررسی قرار گیرد. در نتیجه این فرآیند، آمارهای ارزشمند مختلفی نظیر درصد استفاده از کدینگ های مختلف، توزیع صفحات و وبگاه ها در طبقه بندی های مختلف، سرویس های ارائه شده نظیر جستجو، امنیت و غیره، بدست می آید.

جدول ۱. نمونه ای از ورودی های تجزیه کننده به ازای وبگاه <http://www.whc.ir>

```
<!-- DOCID: 166182 URL: www.whc.ir -->
<head>
<meta http-equiv="Content-Language" content="fa">
<meta http-equiv="Content-Type" content="text/html; charset=utf-8">
<meta name="author" content="xx">
<meta name="copyright" content="2004-2005 by xxx">
<title>Send to friend</title>
<style type="text/css"></style>
<script type="text/javascript"></script>
</head>
<body>
<div id="container">
<form method="POST" action="">
<table>
<tr>
<td>.....</td>
</tr>
<tr>
<td><input type="submit" name="submit">
</td><td>&nbsp;</td></tr>
</table>
<input type="hidden" name="action" value="1">
<input type="hidden" name="newsid" value="1117">
</form>
</div>
</body>
```

۴. نتایج بدست آمده

ورودی سامانه تحلیل گر محتوا، بدین صورت انتخاب شد که به ازای هر وبگاه، حداکثر بیست هزار صفحه مورد بررسی قرار گرفت. علت این امر این است که وبگاه های با تعداد صفحات بالاتر، معمولاً فقط وبگاه های خبری نظیر IRNA، ISNA، IRIB و غیره هستند. در این خصوص، آمار حاکی از آن است که از حدود ۸ میلیون صفحه موجود در دامنه IR، حدود ۶ میلیون آن، مربوط به وبگاه های خبری است. لذا محتوای مناسب به جز خبر، تنها حدود ۲ میلیون صفحه است که این رقم در مقابل اغلب کشورها که بیش از ۱۰۰ میلیون صفحه دارند، رقم بسیار ناچیزی می باشد (به عنوان مثال، دولت الکترونیکی کره جنوبی، شامل بیش از ۱۰۸ میلیون صفحه است).

بواسطه حجم محاسباتی بالای به منظور انجام ارزیابی های مختلف، یک جامعه آماری شامل حدود ۶۰۰ هزار صفحه، در نظر گرفته شد و بررسی های مختلف، روی این مجموعه، انجام شد. در زیر آمار صفحات و وبگاه ها از نظر حجم، عمر و نوع محتوا برای تمام صفحات غیر از اخبار شامل دو میلیون صفحه ارائه شده است.

جدول ۱. نمونه ای از ورودی های تجزیه کننده به ازای وبگاه <http://www.whc.ir>

ورودی سامانه تحلیل گر محتوا، بدین صورت انتخاب شد که به ازای هر وبگاه، حداکثر بیست هزار صفحه مورد بررسی قرار گرفت. علت این امر این است که وب گاه های با تعداد صفحات بالاتر، معمولاً فقط وبگاه های خبری نظیر IRIB ISNA IRNA و غیره هستند. در این خصوص، آمار حاکی از آن است که از حدود 8 میلیون صفحه موجود در دامنه IR. حدود 6 میلیون آن، مربوط به وبگاه های خبری است. لذا محتوای مناسب به جز خبر، تنها حدود 2 میلیون صفحه است که این رقم در مقابل اغلب کشورها که بیش از 100 میلیون صفحه دارند، رقم بسیار ناچیزی می باشد (به عنوان مثال دولت الکترونیکی کره جنوبی شامل بیش از 108 میلیون صفحه است).

بواسطه حجم محاسباتی بالای به منظور انجام ارزیابی های مختلف، یک جامعه آماری شامل حدود 600 هزار صفحه، در نظر گرفته شد و بررسی های مختلف، روی این مجموعه، انجام شد.

در زیر آمار صفحات و وبگاه ها از نظر حجم، عمر و نوع محتوا برای تمام صفحات غیر از اخبار شامل دو میلیون صفحه ارائه شده است.

جدول شماره دو، آمار میانگین حجم و سن (1) و عمق و تعداد صفحات وبگاه ها را نشان می دهد.

عکس

خصوصیات وب ایران ۹۷

۴.۱. آمار وبگاهها

جدول شماره دو، آمار میانگین حجم و سن^۱ و عمق و تعداد صفحات وبگاهها را نشان می دهد.

جدول ۲. خلاصه آمار سایت

تعداد وبگاههای معتبر	۱۰۰۰۰
در داخل کشور ۱۲ درصد وبگاههای	٪۲۷
در خارج کشور ۱۲ درصد وبگاههای	٪۷۳
میانگین تعداد صفحات در هر سایت	۲۳۰
میانگین تعداد صفحات ایستا در هر سایت	۱۲۰
میانگین تعداد صفحات پویا در هر سایت	۱۱۰
میانگین سن مسن ترین صفحه (بر حسب ماه)	۱۳/۶
میانگین سن جوان ترین صفحه (بر حسب ماه)	۶/۸۳
میانگین سن متوسط ترین صفحه (بر حسب ماه)	۸/۵
میانگین عمق ماکزیمم سایت	۴/۵۲
میانگین حجم سایت (بر حسب مگابایت)	۳/۰۶

شکل شماره سه تعداد صفحات وبگاهها را بر حسب درصد سایت نشان می دهد. همانطور که نشان داده شده است بیشتر وبگاهها دارای تعداد صفحات بین ۱۰۰ و ۱۰۰۰ (میانگین ۲۳۰) صفحه می باشند. نمودار دارای توزیع تقریباً خطی می باشد. شکل شماره چهار، توزیع تجمعی شکل شماره سه را نشان می دهد.

شکل شماره پنج، توزیع تجمعی حجم محتوای وبگاهها را نشان می دهد که از این آمار می توان برای مدل سازی محتوای وبگاهها استفاده کرد.

۱. سن صفحه عبارت است از مدت زمان میان ایجاد یا تغییر محتوای صفحه و زمان فعلی

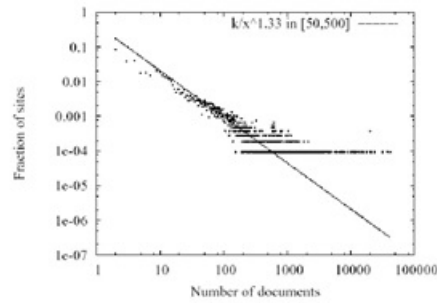
شکل شماره سه تعداد صفحات وبگاهها را بر حسب درصد سایت نشان می دهد. همان طور که نشان داده شده است بیشتر وبگاهها

دارای تعداد صفحات بین 100 و 1000 (میانگین 230) صفحه می باشند. نمودار دارای توزیع تقریباً خطی می باشد. شکل شماره چهار، توزیع تجمعی شکل شماره سه را نشان می دهد.

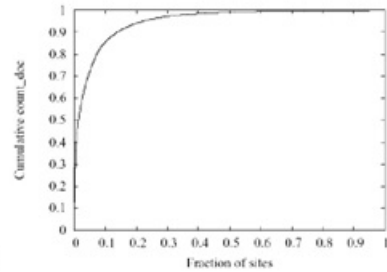
شکل شماره پنج، توزیع تجمعی حجم محتوای وبگاه ها را نشان می دهد که از این آمار می توان برای مدل سازی محتوای وبگاه ها استفاده کرد.

ص: 97

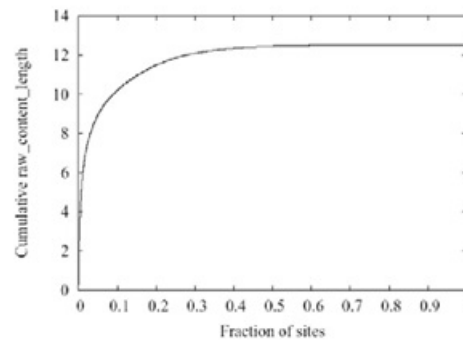
1- سن صفحه عبارت است از مدت زمان میان ایجاد یا تغییر محتوای صفحه و زمان فعلی



شکل ۳. تعداد اسناد در وبگاهها

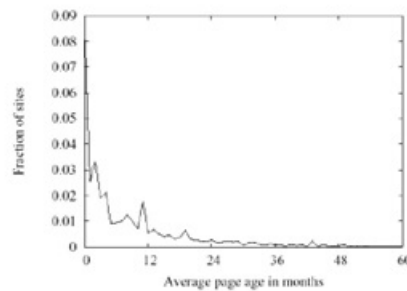


شکل ۴. توزیع تجمعی اسناد در وبگاهها

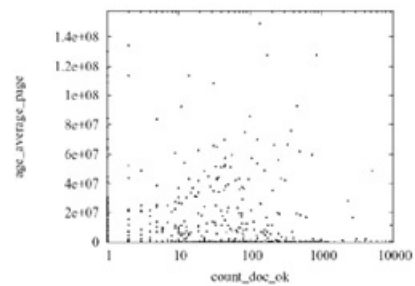


شکل ۵. توزیع تجمعی حجم محتوای وبگاهها

شکل شماره شش نیز میانگین سن صفحات موجود در وبگاهها را نشان می‌دهد. همچنین رابطه بین تعداد اسناد و سن یک سایت در شکل شماره هفت، نشان داده شده است.



شکل ۶. میانگین سن صفحات وبگاهها بر حسب ماه



شکل ۷. تعداد اسناد بر حسب میانگین سن آنها

شکل ۳. تعداد اسناد در وبگاهها

شکل ۴. توزیع تجمعی اسناد در وبگاهها

شکل ۵. توزیع تجمعی حجم محتوای وبگاهها

شکل شماره شش نیز میانگین سن صفحات موجود در وبگاه ها را نشان می دهد. همچنین رابطه بین تعداد اسناد و سن یک سایت در شکل شماره هفت، نشان داده شده است.

شکل 6. میانگین سن صفحات وب گاه ها بر حسب ماه

شکل 7. تعداد اسناد بر حسب میانگین سن آن ها

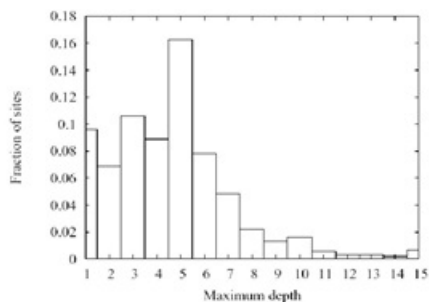
ص: 98

شکل شماره هشت، توزیع عمق وب گاه ها را نشان می دهد. درصد قابل توجهی از وب گاه ها دارای عمق 5 (23%) می باشند. شایان ذکر است که با توجه به شاخص های بدست آمده از وب گاه ها نظیر عمق، تعداد صفحات، حجم محتوا و سن صفحات، به راحتی می توان یک سایت را مدل سازی کرد. با داشتن یک مدل مناسب از وبگاه ها می توان الگوریتم های مختلف را به راحتی تحت آزمون قرار داد.

عکس

۹۹ خصوصیات وب ایران

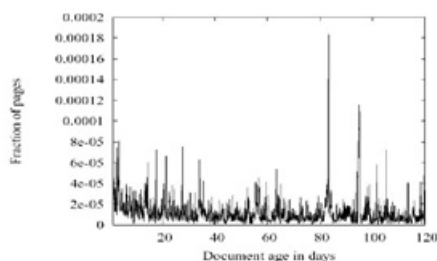
شکل شماره هشت، توزیع عمق وبگاهها را نشان می دهد. درصد قابل توجهی از وبگاهها دارای عمق ۵ (۲۳٪) می باشند. شایان ذکر است که با توجه به شاخص های بدست آمده از وبگاهها نظیر عمق، تعداد صفحات، حجم محتوا و سن صفحات، به راحتی می توان یک سایت را مدل سازی کرد. با داشتن یک مدل مناسب از وبگاهها می توان الگوریتم های مختلف را به راحتی تحت آزمون قرار داد.



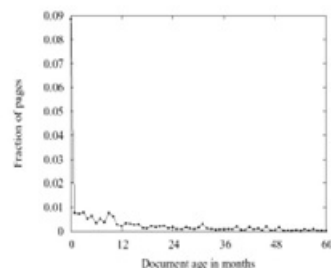
شکل ۸. بیشینه عمق وبگاهها

۴.۲. آمار صفحات

شکل های ۱۰ و ۱۱، به ترتیب، نمایانگر سن درصد صفحات برای حسب روز، ماه و سال است. همانطور که نشان داده شده سن صفحات نسبتاً بالا می باشند. به عبارت دیگر نرخ بروزآوری صفحات، پایین است.



شکل ۹. عمر صفحات بر حسب روز



شکل ۱۰. عمر صفحات بر حسب ماه

شکل 8. بیشینه عمق وب گاه ها

4,2. آمار صفحات

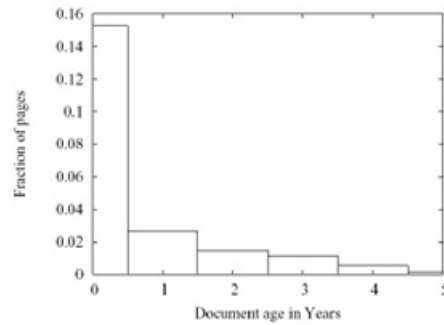
شکل های 9 ، 10 و 11 به ترتیب نمایان گر سن درصد صفحات برای حسب روز، ماه و سال است. همانطور که نشان داده شده سن صفحات نسبتاً بالا می باشند. به عبارت دیگر نرخ بروزآوری صفحات، پایین است.

شکل 9. عمر صفحات بر حسب روز

شکل 10. عمر صفحات بر حسب ماه

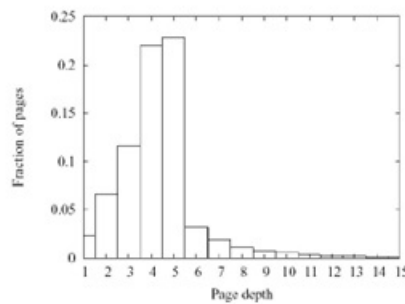
ص: 99

۱۰۰ مدیریت منابع اطلاعاتی وب

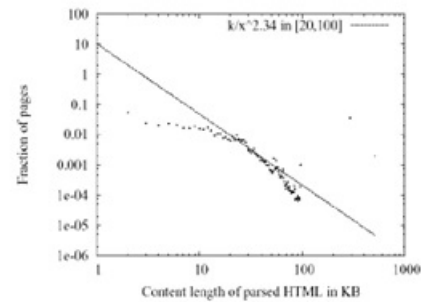


شکل ۱۱. عمر صفحات بر حسب سال

شکل شماره دوازده، عمق صفحات را نشان می‌دهد. با توجه به شکل، اکثر صفحات در عمق ۴ و ۵ قرار دارند. قابل توجه است که عمق صفحات، دارای توزیع پواسون می‌باشد. حجم محتوای صفحات که یک توزیع تقریباً خطی است، در نمودار شماره ۱۳ نشان داده شده است. بیشتر صفحات دارای حجم بین ۱۰ تا ۱۰۰ کیلوبایت را دارا می‌باشند.



شکل ۱۲. عمق صفحات (توزیع پواسون)



شکل ۱۳. حجم محتوای صفحات بر حسب کیلوبایت (تقریباً خطی)

شکل شماره چهارده رابطه بین سن و حجم صفحات را نشان می‌دهد. همچنین رابطه بین عمق با سن و حجم صفحه در شکل‌های ۱۵ و ۱۶ نشان داده شده است. با داشتن شاخص‌های فوق (عمر، حجم، عمق)، به راحتی می‌توان مدل مناسبی از وب ایران بدست آورد.

شکل ۱۱. عمر صفحات بر حسب سال

شکل شماره دوازده، عمق صفحات را نشان می‌دهد. با توجه به شکل، اکثر صفحات در عمق ۴ و ۵ قرار دارند. قابل توجه است که عمق صفحات، دارای توزیع پواسون می‌باشد.

حجم محتوای صفحات که یک توزیع تقریباً خطی است، در نمودار شماره 13 نشان داده شده است. بیشتر صفحات دارای حجم بین 10 تا 100 کیلوبایت را دارا می باشند.

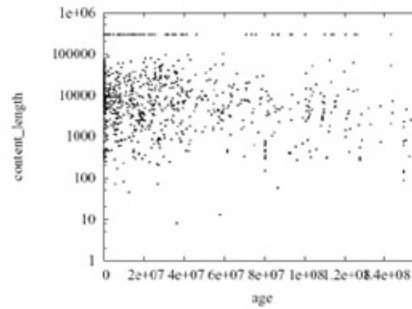
شکل 12. عمق صفحات (توزیع پواسون)

شکل 13. حجم محتوای صفحات بر حسب کیلوبایت (تقریباً خطی)

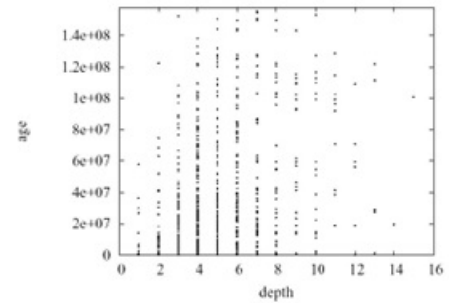
شکل شماره چهارده رابطه بین سن و حجم صفحات را نشان می دهد. همچنین رابطه بین عمق با سن و حجم صفحه در شکل های 15 و 16 نشان داده شده است. با داشتن شاخص های فوق (عمر، حجم، عمق)، به راحتی می توان مدل مناسبی از وب ایران بدست آورد.

ص: 100

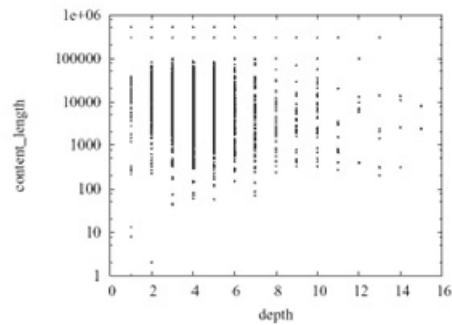
۱۰۱ خصوصیات وب ایران



شکل ۱۴. سن صفحات بر حسب حجم محتوای صفحات



شکل ۱۵. سن صفحات بر حسب عمق



شکل ۱۶. عمق صفحات بر حسب حجم آنها

۴.۳. آمار نوع صفحات

جدول ۳، ۴، ۵، ۶، ۷، ۸ و ۹ به ترتیب، آمارهای مختلفی از صفحات، اعم از ایستا، پویا، صوتی، ویدیویی، تصویری و اجرایی را نشان می‌دهد. طبق جدول ۵، حدود ۵۰ درصد از صفحات، پویا هستند. در مجموع حدود ۸ میلیون آدرس صفحه ایستا و ۱۷ میلیون آدرس صفحه پویا کشف شده است. از بین صفحات پویا، بیشترین درصد را انواع PHP و ASP تشکیل داده‌اند.

جدول ۳. آمار انواع صفحات

تعداد صفحات (بر حسب میلیون)	۲
درصد صفحات ایستا	۵۰/۹۴
درصد صفحات پویا	۴۹/۰۶
درصد صفحات دوتخته‌ای (Duplicate)	۱

شکل ۱۴. سن صفحات بر حسب حجم محتوای صفحات

شکل ۱۵. سن صفحات بر حسب عمق

شکل ۱۶. عمق صفحات بر حسب حجم آنها

4.3. آمار نوع صفحات

جدول 3، 4، 5، 6، 7، 8 و 9 به ترتیب، آمارهای مختلفی از صفحات، اعم از ایستا، پویا، صوتی، ویدیویی، تصویری و اجرایی را نشان می دهد. طبق جدول 5، حدود 50 درصد از صفحات، پویا هستند. در مجموع حدود 8 میلیون آدرس صفحه ایستا و 17 میلیون آدرس صفحه پویا کشف شده است. از بین صفحات پویا، بیشترین درصد را انواع PHP و ASP تشکیل داده اند.

جدول 3. آمار انواع صفحات

ص: 101

جدول ۴. آمار نوع صفحات پویای کشف شده

Filename extension	Number of links found	Percent
php	8,643,428	49.45%
asp	8,429,913	48.22%
jsp	158,477	0.91%
cfm	108,619	0.62%
cgi	81,769	0.47%
shtml	45,847	0.26%
pl	7,507	0.04%
dll	3,130	0.02%
fcgi	1,046	0.01%

جدول ۵. آمار درصد صفحات ایستای کشف شده

Filename extension	Number of links found	Percent
html	7,807,297	95.84%
pdf	220,004	2.70%
xml	63,679	0.78%
doc	15,493	0.19%
rdf	13,951	0.17%
ppt	5,731	0.07%
mso	5,715	0.07%
txt	5,335	0.07%
xls	2,763	0.03%
ps	2,328	0.03%
tex	1,641	0.02%
dvi	1,379	0.02%

جدول ۶. آمار اسناد صوتی کشف شده

Filename extension	Number of links found	Percent
mp3	27,839	44.14%
ram	16,436	26.06%
wma	16,216	25.71%
mid	1,744	2.77%
wav	513	0.81%
pls	173	0.27%
asf	78	0.12%
au	64	0.10%

جدول ۷. آمار ویدیوی کشف شده

Filename extension	Number of links found	Percent
swf	324,678	91.51%
wmv	26,666	7.52%
mpg	2,181	0.61%
avi	1,179	0.33%
mov	99	0.03%
qt	6	0%

جدول ۴. آمار نوع صفحات پویای کشف شده

جدول ۵. آمار درصد صفحات ایستای کشف شده

جدول ۶. آمار اسناد صوتی کشف شده

جدول 7. آمار ویدیوی کشف شده

ص: 102

خصوصیات وب ایران ۱۰۳

جدول ۸. آمار فایل‌های تصویری کشف‌شده

Filename extension	Number of links found	Percent
gif	11,441,330	81.13%
jpg	2,310,009	16.38%
png	300,407	2.13%
ico	38,371	0.27%
bmp	11,536	0.08%
wmf	694	0%
tiff	198	0%
img	2	0%
pbm	1	0%

جدول ۹. آمار نرم‌افزارهای کشف‌شده

Filename extension	Number of links found	Percent
exe	40,344	97.11%
iso	1,066	2.57%
rpm	72	0.17%
patch	42	0.10%
diff	9	0.02%
pdb	9	0.02%
deb	4	0.01%

۵. نتیجه‌گیری و کارهای آینده

هدف این مقاله، بررسی ویژگی‌های وب فارسی است. بر این اساس، سامانه نیمه خودکاری جهت انجام ارزیابی‌های مورد نظر، طراحی و پیاده‌سازی شد. از خصوصیات این سامانه خودکار بودن آن می‌باشد که در زمان‌های مختلف می‌توان آنرا اجرا و آمارهای مورد نظر را استخراج کرد. بررسی‌ها بر روی حدود ۱۱ هزار سایت با پسوند IR و مشتمل بر دو میلیون صفحه انجام شد. این مجموعه، اغلب سازمان‌های دولتی و غیر دولتی را پوشش می‌دهد. هدف اصلی پژوهش، استخراج و تعیین ویژگی‌های وب اعم از خصوصیات ساختاری نظیر ویژگی‌های گراف وب و نیز خصوصیات محتوایی وب فارسی است. شاخص‌های استخراج شده شامل توزیع محتوای وب‌وبگاه‌ها، حجم محتوای فارسی، نوع محتوا، سرویس‌های ارائه شده (جستجو، RSS و غیره) نرخ بروزآوری محتوا، توزیع مکانی وبگاه‌ها در داخل و خارج کشور و غیره می‌باشد. با انجام این پروژه بهتر می‌توان راهبردهای آینده مربوط به ICT را تعیین و تبیین کرد. از نتایج بارز این آمار عبارتند از: کم بودن حجم محتوای فارسی وب در مقایسه با سایر کشورها (نسبت ۱۰٪)، نیاز به یک مرکز داده اینترنتی در کشور (۷۳٪) سایت‌ها در خارج از کشور هستند، نیاز به موتورهای جستجوی بومی (محتوا بیشتر از کدینگ‌های عربی استفاده کرده است)، کم بودن سرویس‌های دولت الکترونیک (جستجو، امنیت و غیره). علاوه بر موارد فوق خروجی‌هایی مانند مجموعه لغات فارسی موجود در وب برای استفاده در خطایاب‌ها و موتورهای جستجو بدست آمده است. در ادامه این کار، در نظر است تا با مطالعه تطبیقی وب فارسی در طی دوره‌های مختلف زمانی،

جدول ۸. آمار فایل‌های تصویری کشف شده

جدول ۹. آمار نرم‌افزارهای کشف شده

۵. نتیجه‌گیری و کارهای آینده

هدف این مقاله، بررسی ویژگی های وب فارسی است. بر این اساس، سامانه نیمه خودکاری جهت انجام ارزیابی های مورد نظر طراحی و پیاده سازی شد. از خصوصیات این سامانه خودکار بودن آن می باشد که در زمان های مختلف می توان آن را اجرا و آمارهای مورد نظر را استخراج کرد. بررسی ها بر روی حدود 11 هزار سایت با پسوند IR و مشتمل بر دو میلیون صفحه انجام شد. این مجموعه اغلب سازمان های دولتی و غیر دولتی را پوشش می دهد. هدف اصلی پژوهش، استخراج و تعیین ویژگی های وب اعم از خصوصیات ساختاری نظیر ویژگی های گراف وب و نیز خصوصیات محتوایی وب فارسی است. شاخص های استخراج شده شامل توزیع محتوای وب و وب گاه ها، حجم محتوای فارسی، نوع محتوا، سرویس های ارائه شده (جستجو، RSS و غیره) نرخ بروزآوری محتوا، توزیع مکانی وبگاه ها در داخل و خارج کشور و غیره می باشد. با انجام این پروژه بهتر می توان راهبردهای آینده مربوط به ICT را تعیین و تبیین کرد. از نتایج بارز این آمار عبارتند از: کم بودن حجم محتوای فارسی وب در مقایسه با سایر کشورها (نسبت 10%)، نیاز به یک مرکز داده اینترنتی در کشور (73%)، کم بودن سایت ها در خارج از کشور هستند، نیاز به موتورهای جستجوی بومی (محتوا بیشتر از کدینگ های عربی استفاده کرده است)، کم بودن سرویس های دولت الکترونیک (جستجو، امنیت و غیره). علاوه بر موارد فوق خروجی هایی مانند مجموعه لغات فارسی موجود در وب برای استفاده در خطایاب ها و موتورهای جستجو بدست آمده است.

در ادامه این کار، در نظر است تا با مطالعه تطبیقی وب فارسی در طی دوره های مختلف زمانی،

نرخ تغییرات شاخص های مختلف تعیین گردد. این اطلاعات، جهت گیری تغییرات وب فارسی را نشان خواهد داد و کمک شایانی به تصمیم گیری های کلان حوزه فناوری اطلاعات خواهد کرد. در این عین حال، در نظر است تا با بررسی های بیشتر در رابطه با محتوای صفحات، اطلاعات بیشتری نظیر متوسط طول حروف در کلمات فارسی، متوسط تعداد کلمات در اسناد، مشخصات گراف (مانند توزیع درجه های ورودی و خروجی) برای زبان فارسی و آمارهای متنوع دیگر، تعیین گردد.

منابع

اشاره

<http://www.worldwidewebsite.com/>, January 2013

Krootkaew C., A. Vongpakaymas, A. Jeawpoung 2002. Services E-readiness Explorer (SEE): Automatic ..monitoring tool for thailand e-government project in proceeding of EurAsia-ICT, Shiraz, Iran, Oct

Baeza-Yates R., C. Castillo and V. LÓpez, 2005. Characteristics of the Web of Spain. Journal of) .Cybernetics, 9(1

Baeza-Yates, R. and F. Lalanne.2004.Characteristics of the korean web. Technical report, Korea Chile IT .Cooperation Center ITCC, 2004

Bordino, I. and D. Donato.2009 .Dynamic characterization of a large Web graph, Dynamic characterization .of a large Web graph. Proceedings of the WebSci'09: Society On-Line, pp. 198-202

Broder A. Z., S.R. Kumar ,F. Maghoul ,P., Raghavan,S. Rajagopalan ,R. Stata , A. Tomkins and J.L.Wiener ..2000.Graph structure in the web. Computer Networks, 33(1): 309-320

Gabriel Tolosa, G., F. Bordignon, R. Baeza-Yates, C. Castillo .2007.Characterization of the argentinian) .web. International Journal of Scientometrics, Informetrics and Bibliometrics, 7(1

Gomes, D. and M.J. Silva.2003. A characterization of the Portuguese Web.In Proceedings of 3rd ECDL .Workshop on Web Archives, Trondheim, Norway

Keyhanipour A.H., A.M. Zare Bidoki, M. Mahmoudi and M. Azadnia.2007.Evaluation of Iran's web content from e-government perspective. Proceedings of the 12th International CSI Computer Conference, .pp. 2081-2086, Tehran, Iran

Rauber, A., O. Aschenbrenner, O. Witvoet, R.M. Bruckner, and M. Kaiser.2002.Uncovering information) .hidden in Web archives. D-Lib Magazine, 8(12

Thelwall, M. and D. Wilkinson.2003.Graph structure in three national–academic webs: Power laws with anomalies. Journal of the American Society for Information Science and

ص: 104

.Technology, 54(8): 706-712

Thelwall, M. and A. Zuccala. 2008. A university-centred European Union link analysis. *Scientometrics*,
.75(3): 407-420

Shestakov, D.2011.Sampling the national deep web, database and expert systems applications. Lecture
.Notes in Computer Science, 6860: 331-340

ص: 105

مقاله حاضر گزارشی است که توسط پژوهشگران مؤسسه اینترنت آکسفورد برای کنفرانس بین المللی حفاظت اینترنت تهیه شده است. هدف از آن ایجاد انگیزه بحث و تبادل نظر بیشتر مابین آرشئویست ها و پژوهشگران وب در مورد روش های آرشئو وب است در بخش اول مقاله نمای کلی 4 سناریو احتمالی از آرشئوهای وب شامل نیروانا، آپو کالپس انفرادی و آرشئوهای غبار آلوده را معرفی می کند. سپس به توصیف انواع مختلفی از پژوهش هایی که در مورد وب پویا در حال انجام هستند، می پردازد. و در بخش آخر چالش های فعلی و پیش روی آرشئوهای وب را پوشش می دهد و در نهایت به ارائه پیشنهاد های برای رفع چالش ها و ارائه راه حل های میان مدت و دراز مدت برای مقابله با تغییرات احتمالی را ارائه می کند.

*آینده آرشیو وب (1)

نوشته: اریک تی. مهیر (2)، آرتور توماس (3)، رالف شرودر (4) | مؤسسه اینترنت آکسفورد (5)

ترجمه: رضا خانی پور (6)، محبوبه قربانی (7)

خلاصه اجرایی

این گزارش، توسط پژوهشگران مؤسسه اینترنت آکسفورد برای کنسرسیوم بین المللی حفاظت اینترنت (8) به نگارش در آمده است. هدف از آن ایجاد انگیزه بحث و تبادل نظر بیشتر مابین آرشیویست ها و پژوهشگران وب در مورد روش های آرشیو وب است که می تواند مورد استفاده پژوهشگران قرار گیرد.

بخش اول. نمای کلی از چهار سناریوی احتمالی برای آینده:

اشاره

- نیروانا (9): مکانی که آرشیوهای وب توسط بسیاری از گروه ها به طور گسترده استفاده، استاندارد سازی، و قابل مرور می شود و رابط کاربر قوی برای دسترسی دارند.

ص: 107

Web archives: the future(s) -1

Eric T. Meyer -2

Arthur Thomas -3

Ralph Schroeder -4

Oxford Internet Institute -5

6- عضو هیئت علمی سازمان اسناد و کتابخانه ملی ج. ا. ا.

7- دانشجوی دکترای علم اطلاعات و دانش شناسی واحد علوم و تحقیقات تهران دانشگاه آزاد اسلامی

8- (the International Internet Preservation Consortium (IIPC

9- Nirvana

- آپوکالیپس (1): آرشیوها تجزیه می شوند و استاندارد نشده باقی می مانند، به سختی بازیابی شده و در دسترس قرار می گیرند و در نتیجه مفید نیست و به سختی استفاده می شوند.

- انفرادی: در این سناریو، آرشیوها مانند یک هوش هم بسته منفرد که قادر است بین اشیای رقومی و افراد ارتباط برقرار کند، غیر ضروری می شوند.

- آرشیوهای غبارآلود: در این سناریو، جامعه آرشیوی وب هرگز پاسخگوی سؤال (که چی؟) نیست آرشیوهای وب به طور گسترده، مورد استفاده قرار نمی گیرد، بلکه در حال جمع آوری غبارهای و رقومی هستند.

این سناریوها ما را به تفکر درباره تعامل بین آرشیوها پژوهشگران و محققان در روش های گوناگون قادر می سازد.

بخش 2. انواع مختلف پژوهش هایی را توصیف می کند که در مورد وب پویا در حال انجام هستند. فنی که حالا به طور گسترده ای، بیشتر از استفاده از آرشیوهای وب مورد بهره برداری قرار می گیرد.

مقصود این است که استفاده از وب پویا می تواند الهام بخش تفکر درباره کاربردهای بالقوه آرشیوهای وب باشد. این کاربردها عبارت اند از:

مصورسازی: که از طریق آن پیوندها، نه تنها بین وبگاه ها، بلکه میان انواع مختلف اطلاعات نیز

می توانند ایجاد شوند، همانگونه که سازمان ها و دیدگاه های کلی آرشیوها را توانمند می سازد.

سنجش های دگرساز (2): دانشمندان حوزه علم سنجی شروع به کسب داده ها از منابع مجزا، از طریق تحلیل های استنادی نموده اند. به طور مثال: بلاگ های محققان و پیوندهای بین این بلاگ ها.

فنون دیگری، همچون نقشه برداری محتوای ایجاد شده توسط کاربر و تحلیل های شبکه اجتماعی در اینجا ارائه شده است.

بخش 3. چالش های فعلی و پیش رو را پوشش می دهد. نخستین قسمت از این بخش روش های تغییر وب، پیشنهادهایی در این مورد و راه حل های میان مدت و دراز مدت برای مقابله آرشیوها با این تغییرات را توصیف می نماید. این بخش با پیشنهادهایی درباره گام های پیش رو پایان می یابد.

مقدمه

این گزارش، به صورت پیش نویس در می 2012 نگارش یافت و در مجمع عمومی آی آی. پی. سی. (3) در

ص: 108

(<http://altmetrics.org/about>)

3- IIPC (آی آپی سی) کنسرسیوم بین المللی حفاظت اینترنت (<http://www.netpreserve.org>) که این کار را پایه گذاری و زیربنایی برای بحث های بیشتر ارائه نمود که با ارائه در نشست مجمع عمومی IIPC در سال 2011 و برگزاری کارگاهی که ویرایشی نهایی از این گزارش خواهد بود آغاز می شود.

اين گزارش، در يك نشست جامع، تلخيص و در يك كارگاه نيز ارائه شده است. همچنين، از طريق پست الكترونيك براي جامعه محققان اينترنت و جامعه كتابداري و اطلاع رسانی ارسال شده است. هدف اين بود كه پيش نويس گزارش منجر به تهيج و ترغيب شود، افكار را برانگيزد. براي اينكه آرشيويست هاي وب و پژوهشگران را به خروج از سكون وادار كند. افراد مشابه و ديگران را به فعاليت وادار كند. براي ايجاد تغيير، پيش از اين براي بحث و تبادل نظر ايجاد انگيزه شده است؛ كه فعاليت را به سوي تغيير ملموس ترغيب مي كند.

چرا تغيير ضروري است؟ وقتی كه IIPC توجه ما را به اجراء اين پروژه جلب كرد، اين ضرورت كه جامعه آرشيويست وب و IPC به طور ويژه، بايد روش هاي نوين براي تشويق کاربران و کاربردهاي جديد آرشيوهاي وب و مدل هاي جديد آرشيوسازی وب و سبك هاي جديد تعامل با محققان بررسی نماید، احساس می شد.

اين مباحث، پيش تر در قالب دو گزارش به وسيله (2) JISC تدوين شده بود كه بر روي شرايط فعلي هنر آرشيوهاي وب (دو گرتي و ديگران، 2010) (3) و فرصت هاي سرمايه گذاري جديد متمرکز شده بود.

بعضی از نتايج اين دو مقاله در ادامه آمده است، اما موضوعات عمومي سرتاسر اثر اين بود كه هنوز فاصله و شكافي بين جامعه پژوهشگران بالقوه اي كه دليل خوبي براي استفاده، تحليل و اشتراك آرشيو وب دارند و جامعه واقعي (به طور كلي هنوز كوچك) پژوهشگران كه در حال حاضر فعاليت مي كنند وجود دارد (دو گرتي و ديگران، 2010، ص.5). تجربه كار ما در اين گزارش و صحبت با اعضاي آي.آي.بي.سي. و جامعه پژوهشگران اينترنتي اندكي باعث تغيير نظر ما در مورد اين موضوع شده است، مسلماً ما از اينكه هميشه کاربردهاي آرشيوهاي وب به خوبي مورد بحث قرار گيرد و جامعه پژوهشي هنوز به صورت معني داري به آن دست نيافته اند، اطمينان بيشتري يافته ايم. اين گزارش، به نوبه خود اين [ديدگاه] را كمی تغيير خواهد داد، اما اگر پاره اي از پيشنهادهاي ارائه شده در آن به طور جدی توسط جوامع مرتبط، به كار گرفته شود احتمال اينكه در آينده آرشيوهاي منابع اينترنتي براي پژوهشگران اهميت بيشتري يابد، افزايش خواهد يافت.

اين گزارش، در ابتدا به منظور اينكه افكار نظري را به آينده احتمالي وب معطوف كند، به صورت يك تمرين كه ما را وارد به تفكر درباره آن چه كه براي انجام در حال حاضر به آن نياز داريم و اطمينان از اينكه بتوانيم به صورت مطمئن و پر ثمر از آرشيو وب در آينده استفاده كنيم، شكل يافته است. سپس، ما بر روي روش هاي ابزاري مورد استفاده براي پژوهش وب پويا تمرکز می کنیم، همانند يك اشاره گر بر روي انواع چيزهائي كه می تواند براي كمك در فهم وب آرشيو، گسترش يابد. در ادامه، ما بر روي يك سري از موضوعات و سؤالاتي كه پژوهشگران می خواهند يا ممكن است بخواهند براي استفاده از وب آرشيو

ص: 109

the Hague - 1

2 - ISC كميته سيستم هاي اطلاعاتي مشترك (<http://www.jisc.ac.uk/>) كه بخش پژوهش حوزه اطلاعات و ارتباطات و توسعه زيربنايي آموزش و پژوهش در بریتانیا را پایه گذاري نمود.

Dougherty - 3

به آن‌ها بپردازند تمرکز می‌کنیم. در پایان این بخش، برخی از چالش‌های انفرادی، سازمانی و سازمان‌های بین‌المللی که می‌توانند برای افزایش توانایی ما برای توضیح این موضوعات و پاسخ به این سؤالات، مورد توجه قرار گیرند را تبیین می‌نماییم.

ما گزارش مذکور را با نتایج بر اساس آن چه از این تجربه آموخته‌ایم، به پایان می‌بریم.

ساختن آینده

(بهترین راه برای پیش‌بینی آینده، ساختن آن است) (کای، 1995) (1)

برای آغاز بحث از آینده‌شناسی بهره می‌بریم. هدف پیش‌بینی آینده نیست. که اشتباه خواهد بود. در واقع، ما کاملاً در مورد تلاش‌هایی که برای پیش‌بینی آینده انجام شده است، اطمینان نداریم. ساختن سناریوها و تلاش‌های دیگری که برای پیش‌بینی آینده (ادعاهایی درباره آینده) طراحی شده‌اند، معمولاً در دانش، نهفته می‌مانند، به طوری که بیشتر چنین تلاش‌هایی، هرگز جدی گرفته نخواهد شد.

با وجود این، حداقل یک نوع از آینده‌شناسی که از نظر ما مناسب است وجود دارد. [و] آن زمانی است که هدف از تجربه، پیش‌بینی آینده نباشد؛ بلکه الهام بخشی باشد برای افرادی که مسئول ساختن نظام‌هایی هستند که زیربنای آینده خواهند بود. به منظور تفکر درباره نتایج تصمیمات فعلی در دوره‌هایی با تأثیرات احتمالاً طولانی مدت. آی.پی.سی. متشکل از بسیاری از افرادی است که در حال حاضر عهده دار توسعه و گسترش نظام‌ها، ابزارها، استانداردها، و تفاهم‌نامه‌هایی برای حفاظت از محتوای اینترنت با نیم‌نگاهی به سوی مفید ساختن آن برای فهم جامعه‌ای که در آن زندگی می‌کنیم، هستند.

در گسترش نظام‌های رایانه‌ای، «انتخاب‌های معماری» بسیاری وجود دارد. (کلینگ، مک‌کیم و کینگ 2003) (2)، در طول مسیر نقاطی که تصمیمات در آن جا اتخاذ می‌شوند، همان‌هایی هستند که یکی در دو راهی‌های طول مسیر به غیر از انتخاب‌های دیگر گزیده می‌شوند. شواهدی وجود دارد از اینکه جامعه آرشیوی وب اهمیت نقاط انتخاب در حال حاضر و در آینده نزدیک را می‌پذیرد. (پذیرفته است). یک سری از نقاط انتخاب با پیشنهادهایی که نوید تغییر زلزله‌گونه‌ای در آرشیو وب می‌دهند، مرتبط هستند: گذر از دسترسی به وبگاه‌ها و صفحه‌های انفرادی به سوی ساختن و استفاده از یک مجموعه همانند یک مجموعه در مقایسه با ایجاد دسترسی به بخش‌هایی از یک مجموعه. به چه تصمیماتی برای افزایش احتمال اینکه آرشیو - در - جعبه، مفید، قابل استفاده، ادامه‌پذیر و تاثیرگذار باشد، نیاز خواهد بود؟ ما دوباره به ایده آرشیو - در یک جعبه در این مقاله، موازی با چالش‌های دیگری که انتخابها را می‌طلبد، باز خواهیم گشت.

در ادامه، هدف این تمرین این است که جامعه آرشیوی وب کدام راه را برای آن انتخاب‌ها که بر آینده تاثیر خواهند گذاشت و برای پیشنهاد گام‌ها و انتخاب‌هایی که آینده را در یک جهت یا جهت‌گیری هدایت [ترغیب] می‌کند، درخواست می‌کند.

ص: 110

Kay - 1

Kling, McKim, King - 2

اشاره

ما می توانیم آینده های بسیار مختلفی را برای آرشیوهای وب و کاربردهای شان تصور نماییم، برای فرض بحث، ما چهار سناریوی بالقوه را که می توانند در دهه یا دو دهه آینده اجرا شوند، بررسی الزامات شان و راه های پیشنهادی که جامعه آرشیوی وب باید با آن ها هماهنگ شوند، ارائه داده ایم:

در ادامه این متن، به طور ویژه تعدادی از عناصر که این سناریوها را تکمیل خواهند کرد، نشان می دهیم، چالش هایی را که در سر راه به کارگیری آن ها وجود دارد شناسایی می نمائیم، و مثال هایی از گستره متنوع ابزارهای توسعه وب پویا را نشان می دهیم. اگر آن ها برای داده های تاریخی به کار گرفته شوند، بین بهترین و بدترین طرح می توانند تفاوت ایجاد نمایند.

سناریوی نیروانا

در بهترین شرایط، آرشیوهای وب، مستحکم، استاندارد شده و به صورت ایمن، حفاظت شده خواهند بود. در حالی که، در زمان مشابه به صورت باز، انعطاف پذیر و گسترده، و به عنوان بخشی از ابزار استاندارد پژوهشی در علم اینترنت، علوم سیاسی، اقتصادی، جامعه شناسی، تاریخ معاصر (و در آینده، تاریخ اواخر قرن بیستم و اوایل قرن بیست و یکم) روزنامه نگاری، زبان شناسی، ارتباطات، تجارت، مطالعات رسانه و دیگر رشته ها استفاده می شوند.

علاوه بر دانشگاه، آرشیوهای وب برای عموم، حکومت ها، واحدهای سیاسی، و گروه های مشاوران

(متفکران) و سازمان های غیر دولتی مفید و کاربردی خواهد بود. متأسفانه، از خیلی از جنبه های بسیاری، کم احتمال ترین سناریوست، به این سبب که به نظر می رسد در حال حاضر، برای امکان پذیر ساختن ورود آن به جامعه آرشیوی وب به تلاش گسترده تر و منابع وسیع تری نیاز خواهد بود. اگر چه، ممکن است برای اینکه آن را در ذهن خود به صورت مطلوب نگهداری کنیم، مفید به نظر برسد، همان گونه که ایجاد توازن را بین آن چه که می توانیم و آن چه باید باشد بررسی می کنیم.

به منظور ارائه این سناریو، حتی در طرح کلی، نیاز است که اتفاق هایی رخ دهد (در بخش بعدی، مثال هایی از وب پویا در زمان وقوع، ارائه خواهیم کرد) این موارد عبارت اند از:

- توسعه قوی تر و مؤثرتر ابزارهای جست و جوی متن، تحلیل و استخراج اطلاعات، مصورسازی، وبلاگ نویسی اجتماعی، تحلیل های طولی و تحلیل های نظری.

- توسعه روش های بهتر برای اینکه کاربران، گشتالت مجموعه های چندتایی و انفرادی را درک کنند. در حالی که محتوای متنی می تواند جست و جو شود، نیاز به فراداده غنی برای پشتیبانی خطوط کلی محتوا و یا روش های جدید سازماندهی آن می باشد.

افراد به طور ویژه در تشخیص الگوهای بصری توانمند هستند، و به نظر می رسد ابزارهای گرافیکی یکی از بهترین راه های کسب آن می

باشند. می‌توانیم ایجاد فضای مجازی را تصور کنیم که به [fly through 1](#) های سه بعدی و سایر روش‌های فضایی حسی امکان سازماندهی محتوا را می‌دهد. در بهترین

ص: 111

1 - Fly through: نوعی شبیه‌سازی کامیوتری است که در فضاهای مجازی به کاربران اجازه مشاهده مدل‌های مجازی سایت‌ها را می‌دهد. (دسترسی در: <http://dictionary.reference.com/browse/fly-through>)

حالت، محیط های مجازی کاملاً همه جانبه از نوع غار (1) (شرودر، 2011) می توانند همکاری گروه های توزیع شده فضایی (مکانی) افراد را پشتیبانی کرده و امکان اشتراک مؤثر و تعامل اجتماعی را ایجاد کنند. از این دیدگاه، تمامیت آرشیوهای وب می تواند به صورت یک کلیت بزرگ در نظر گرفته شود، که افراد می توانند در آن به صورت منفرد یا گروهی سیر نمایند. وب فعلی (و به همین ترتیب آرشیوهای وب) درک ناچیزی از سازماندهی فضایی ارائه می دهند و نشانه های مناسب فضایی یا قابلیت های دیگر کمک به مرور و کاوش را محدود (سلب) می کنند. بر خلاف کتابخانه فیزیکی، اسناد رقومی در محیط جدا سازی شده از یکدیگر نگهداری می شوند، در مکانی که اسناد در دنیایی سه بعدی سازماندهی می شوند تا امکان مرور آن ها به واسطه گردش (جست و جو) در محیط را میسر ساخته و نشانگرهایی همچون همجواری فضایی (مکانی) به کشف آن ها کمک می کنند.

- در حالی که ابزارهای وب نوشت نویسی اجتماعی پدید می آیند، آرشیو های وب، فاقد سایر ابزار های همکاری هستند، مانند موتورهای پیشنهادی (مخزن پیوسته آمازون مثال عالی از این امکان است).

- به طور فزاینده ما نیازمند محتوای آماده شده توسط کاربران (وب 2) در مقیاسی بزرگ (فیس بوک (2)) هستیم. اما ساختاری حیرت انگیز بر اساس چنین محتوای سازمان نیافته نامتجانس و ذاتاً سازمان نیافته ای، در این مقیاس با مشکل مواجه است. انجام آن به وسیله ماشین از نظر فناوری بسیار چالش برانگیز است، بخشیدن غنای معنایی، به عنوان جایگزینی برای پشتیبانی رهیافت انباشت منبع، به کاربران آرشیو اجازه می دهد که محتوا را سازماندهی کنند این شکل نهایی وب نوشت نویسی (حاشیه نگاری) اجتماعی است که کاربران نه فقط داده بلکه فراداده را نیز تولید می کنند (گازان، 2008) (3).

در نیروانا، انتخاب های امروز توسط پژوهشگران آینده ستایش خواهد شد، کسانی که به اتکای اطلاعات و شواهد، تلاش های انسانی انجام شده در اینترنت را، تجسم داده، حفاظت کرده، و افزایش داده اند تا به همه روش های فنون قدرتمند پژوهش دست یابند و توانمند شوند.

سناریوی آپوکالیپس (آخر الزمان)

در بدترین حالت، تغییرات روزافزون اینترنت برای بسط و توسعه فناوری های جدید (HTML5)، محتوای قابل اجرا، ویدئوی جاسازی شده و اشیای تعاملی، وبگاه های مبتنی بر پایگاه داده، نرم افزارهای موبایل غیر وابسته به HTML.TTP و مانند آن) با سرعتی گیج کننده، ادامه خواهند یافت، و ابزارهای آرشیو سازی وب در حفظ همگامی با آن ها شکست خواهند خورد، و بیش از پیش عقب می افتند. حتی اگر فناوری های آرشیوسازی وب بتوانند سرعت را حفظ کنند، تغییرات دائمی (پیوسته) همه اشکال مطرح شده، چالشی حل نشدنی را ایجاد می نماید. در این طرح به طور صادقانه ای، فقط اندکی از محتوای واقعی ذخیره می شود، و حتی اگر ذخیره شود، افزودنی های اختصاصی برای مشاهده آن حفظ نشده یا قابل حفظ

ص: 112

شدنی نیستند، و مشاهده محتوا غیر ممکن می شود. بیشتر پیشینه های پیوسته در طول دوره ما سرانجام همچون کارت پانچ های دهه 1960 و نوارهای مغناطیسی ریل به ریل (1)، غیرقابل خواندن (استفاده) خواهند شد. علاوه بر مشکل قالب، مشکل حجم زیاد نیز در حال رشد و فزونی است. همان طور که اینترنت به سمت تکامل استفاده از IPv6 (2) حرکت می کند، اشیای قابل نشانه گذاری (به طور روزافزونی شامل، اشیای فیزیکی در «(3) Web of Things» به راستی، به وسیله دستورات متعدد ذخیره سازی - حتی نشانی ها - بیش از ظرفیت موجودمان، در حال ازدیاد حجم هستند (1038)، که به تنهایی اجازه آماده سازی همه محتوا را می دهد. در نتیجه، ما نمی توانیم موارد بیشتری را جست و جو کنیم، و به این دلیل نمایه سازی و فناوری جست و جوی ما ناامیدانه شکست می خورد.

همزمان با توسعه وب معنایی، همه مفهوم محتوا تغییر می کند. محتوا چیزی بیشتر از متن و تصویر نیست، اما در حال حاضر اقلام داده قراردادی و پیوندهای بین آن ها را در بر می گیرد. حتی حالا در سال 2011، جهان داده های پیوند خورده عمومی (4) (که شامل مجموعه هایی همچون data.gov.uk و data.gov است). ده ها میلیون اقلام داده ای (به طور روز افزون قالب (5) RDF) را در بر می گیرد که به وسیله صدها میلیون پیوند قابل ارجاع مجدد، به هم پیوسته اند. چالش آرشیوسازی این سری داده ها از زمان نشانی دهی - آغاز می شود (برای مثال: به وسیله پروژه پرونوم (6) آزمایشگاه های آرشیوهای ملی بریتانیا (7))، اما احتمالاً جهان داده های پیوند خورده سریع تر از اینکه از نظر مجموعه سازی یا تحلیل، قابل مدیریت باشد، رشد خواهند کرد.

در این سناریو، حتی منابع عظیم شرکتی چون گوگل تحت الشعاع مشکلات قرار می گیرد. بنابراین، پاسخ بدیهی برای آرشیوسازی این حجم («بگذارید گوگل کارش را بکند» (8)) چیزی بیشتر از یک راه حل نیست.

اگر انتخاب ها ما را به این مسیر هدایت کند، پژوهشگران آینده می آموزند که گذشته وب را به صورت غیر قابل دسترسی و بدون اعتماد و داستانی که حکایت می شود و مدرک دست دومی از زمان، تصور کنند. امروزه، حجم بزرگی از اطلاعات به مقیاس جهانی ایجاد می شود که ممکن است به صورت نوشته های روی تکه های کاغذ در میلیون ها جعبه کفش ذخیره شده باشد، که همگان، پیشرفت های صورت گرفته در دنیا را همان گونه که در محتوای اینترنت منعکس می شود به نحو مناسبی درک خواهند کرد.

ص: 113

reel to reel -1

2- IPv6 آخرین نسخه پروتکل ارتباطات در اینترنت (دسترسی در <http://en.wikipedia.org/wiki/IPv6>)

3- Web of Things (http://en.wikipedia.org/wiki/Web_of_Things)

4- Linked Data - Connect Distributed Data across the Web, at www.linkeddata.org

5- (Resource description Framework) RDF (چارچوب توصیف منابع)

6- PRONOM

7- <http://www.nationalarchives.gov.uk/PRONOM/Default.aspx>

8- let Google do it

در دنیایی کاملاً متغیر، افراطی ترین سناریو سناریوی انفرادی است که در آن اینترنت همانگونه که ما می شناسیم به سمت پدیده ای کاملاً جدید و احتمالاً در نوع خود هوشمند (کورزویل 2005) (1) تکامل می یابد. در مسیر منحصر به فرد شدن، به سمت توسعه و تبدیل شدن به موجود مجازی پیچیده ای حرکت می کند، که ما فهم کمی از آن داریم و برای آرشو کردن آن، در حال حاضر روشی بهتر از آن چه که با هوش انسانی خود از آرشو می دانیم، وجود ندارد. حتی اکنون، در سال 2011، برای غلبه بر آن آغاز به تشخیص تمایز بین پردازش انسانی و مصنوعی کرده ایم. خدماتی همچون: ریکاپچا (2) (ون آون (3)، مورر (4)، مک میلن (5)، آبراهام (6) و بلام (7) و آمازونز مکانیکال ترک (8) که از موجودیت های انسانی «در یک چرخه» برای حل مشکلاتی که برای ماشین ها سخت است، استفاده می کنند، به ما راه هایی را به سوی دنیایی که در آن هوش انسان و ماشین به صورت دوقلوهای جدایی ناپذیری هستند، نشان می دهند و مرز بین آن ها نامشخص است. در چنین دنیایی، مشخص نیست که آرشوسازی چه مفهومی می تواند داشته باشد. بنابراین، همزمان با حرکت زمان به سمت جلو، گذشته به ناچار و به صورت جبران ناپذیری از دست می رود. این سناریو، ممکن است شبیه به یک داستان علمی به نظر برسد. اگر چه، ارزشمند است که به خاطر داشته باشیم آینده غیر قابل پیش بینی است، حتی اگر تأثیر بر روی شرایط انتخاب در طول مسیر را برای اینکه در یک مسیر یا مسیر دیگری قرار گیرد - مدیریت نماییم. انتخاب هایی که ما انجام می دهیم ممکن است، برای اجرای وظیفه ای جدید، همچون اینترنت هوشمند ناکافی باشد. این بدین معنا نیست که ما باید انتخاب های صحیح را در نظر بگیریم و تلاشی از خود به خرج ندهیم.

سناریوی غبار آلود

متأسفانه این سناریو، احتمالاً در ابتدای امر [این اندیشه را به ذهن متبادر می کند] که آرشوهای وب معادل آرشو رقومی غبار آلود است، اگر چه اغلب به خوبی حفاظت و نگهداری می گردد. به سختی از آن استفاده می شود.

حتی اگر جامعه آرشوسازی، وب، به توسعه استانداردها و تجربیات برای حفظ بخشهای اینترنت ادامه دهد، کاربردهای مؤثر ناچیزی از جامعه پژوهشی سر می زند و جامعه پژوهشی به طور مؤثر از آن استفاده نمی کند.

ممکن است به صفحه های اینترنت به صورت مجزا با استفاده از ابزارهای پیوسته رجوع شود و بعضی از پژوهشگران به ایجاد آرشوهای کوچک برای موضوعات پژوهشی خاص مبادرت ورزند، اما پژوهش اینترنت تنها با اولویت تمرکز بر روی وب پویا ادامه خواهد داشت و در آینده نزدیک توجه اندکی به

ص: 114

Kurzweil -1

2 - reCAPTCHA: برنامه ای برای رقومی کردن کتب، نشریات و ... دسترس —————ی در

(<http://www.google.com/recaptcha/learnmore>)

Von Ahn -3

Maurer -4

Mc Millen -5

Abraham -6

Blum -7

(Amazon's Mechanical Turk(Amazon Corp., Mechanical Turk at www.mturk.com -8

استفاده از وب قبلی برای پژوهش جدی، گسترش خواهد یافت.

این سناریو با سناریوی آخر الزمان متفاوت است. در آن سناریو، فناوری آرشیوسازی وب با تغییرات فناورانه روی اینترنت همگام نخواهد بود. در این سناریو سرعت آرشیوسازی وب هماهنگ با فناوری تحویل وب تنظیم می شود.

اگر چه داده حفاظت شده به صورت اسطوره ای حفاظت شده برای استفاده در آینده نامطمئن باقی می ماند.

در مراحل نگارش این گزارش، مشخص شده است که به جای آرشیوهای وب مرجع، به طور روزافزونی کاربران و پژوهشگران با وب پویا همانند آرشیو برخورد می کنند. وب پویا، به رشد خود ادامه می دهد و برای بیشتر بخش ها، داده هایی که ناپدید می شوند، از نظر بسیاری به صورت یک مشکل ساده دیده می شود و مهم تر از آن، برای بخشی بزرگ تر، حجم عظیم دیگری از داده ها بر روی وب در هر زمان باقی می ماند.

تصویر ما از آرشیو، تصویری از اقلام فیزیکی همچون کاغذها و اسناد ذخیره شده در یک مکان فیزیکی است. با وجود این، این ماهیت وب نیست. پژوهشگران می توانند حجم زیادی از مواد را از منابع مختلف وب پویا برای نیل به اهداف پژوهشی خود ذخیره کنند. ما آرشیوها را به صورت چیزهایی که برای آیندگان حبس شده اند، درک می کنیم. در حالی که وب، به تنهایی پدیده ای در حال رشد و منبعی متنوع از انواع مواد است که به صورت بالقوه مورد توجه پژوهشگران است. به طوری که، آن را نه به عنوان یک آرشیو سنتی، بلکه به سادگی به صورت یک منبع داده می بینند.

این یک سناریوی بدبینانه است، اما به نظر می رسد که وزنی از شواهد را در جانب خود داشته باشد. در راینی هایی با محققان برجسته، به یک فقدان علاقه دیرپای بر پرسش از وب قبلی و درک اینترنت به عنوان یک توسعه تاریخی، رسیدیم. البته استثناهایی وجود دارد که بعداً در این گزارش شرح داده خواهد شد، اما قادر بوده ایم تا گرایش نهفته کار با آرشیوهای وب را کشف کنیم که انتظار ساده ای برای [ظهور] فناوری مناسب برای بیدار شدن آن وجود دارد.

ممکن است انتظار تغییر از گوشه و کنار برای آمادگی ایجاد یک گام تغییری در تصورات پژوهشگران، بر اساس جلوه های استفاده جدید یا فناوری تازه، وجود داشته باشد. اگر این امر رخ ندهد، بیم آن می رود که به جمع آوری گردوغبارهای رقومی ادامه دهیم.

اگر از این سناریو اجتناب شود، نیازمند نوعی جدیدی از آرشیویست هستیم که با پژوهشگران و عموم [کاربران]، در استخراج داده هایی که آن ها از وب پویا نیاز دارند در تعامل باشد و زمانی که این داده ها از وب پویا ناپدید شدند، قادر به ذخیره سازی مجدد آن ها به صورتی که قابل مشاهده و استفاده بوسیله ابزارهای وب پویا باشند، هستند. بسیار فراتر، از زمانی که رقومی سازی اسناد تاریخی، از آرشیوها حجم زیادی از مواد تاریخی قابل دسترسی بر روی وب در دهه گذشته ایجاد کرده است (مهیر 2001 (1))،، تانر 2010 (2))، تانر و دیگران 2011 (3))، آرشیوهای وب به انتقال وب به محفظه ها (جعبه ها) نیاز ندارند، بلکه در عوض نیازمند برگشت به وب در زمانی که دارای محتواست، هستند.

ص: 115

Tanner -2

Deegan -3

در حین این که ما مسیر خود را به سوی آینده آرشیو وب مرور می کنیم، پرسش های متعددی در مورد چگونگی دیدگاه مان نسبت به آینده، می توانیم از خود بپرسیم.

برای مثال، آیا آرشیو آینده، باغ دیوار کشیده شده ای خواهد بود که از آسیب ها محفوظ و در امان است، اما با دسترسی محدود است؟ آیا فضایی کاملاً باز و قابل دسترسی برای تازه واردان خواهد بود؟ آیا مجموعه ای از مخازن است یا یک عرصه ارتباط متقابل منفرد؟ یا اینکه یک شهر ارواح خالی از سکنه باز و مرتبط خواهد بود؟

بخشی از آرشیو ها برای محققان و دانشمندان و بخشی دیگر برای عموم است (1). آیا این دو می توانند در فضایی سایبرنتیک با هم مرتبط باشند - تا اینکه دانشمندان به صورت دائمی، آن چه که عموم [شامل دانشمندان] به آن دسترسی دارند و از آن استفاده می کنند و آن را ایجاد می کنند، پایش نمایند (بنابراین درک از آگاهی جهانی یا دریافت جمعی را افزایش می دهند) در حالی که در زمان مشابه شکل گیری چنین فضایی با این سبک، برای گسترش، مجموعه سازی و تأثیرگذاری و لذت بخش بودن و کاربرد غنی همه دنیا، غنی سازی می شوند؟

حتی در صورتی که منحصر به فرد بودن، با شکست مواجه شود، اینترنت به طور فزاینده ای اجازه ارتباطات در این سبک را که از نظر استفاده برای هوش جهانی خیلی کم قابل فهم می باشد، می دهد (شرودر و مه ی، 2009). در هوش جهانی - با خوراک ها و پیوندهای پویا - هوش های انفرادی از طریق ابزارهای ورودی و خروجی به یکدیگر مرتبط هستند. چگونه آرشیوهای وب ماهیت ارتباطات درونی هوش جهانی را به صورت غیر از اسناد به ظاهر جدا از هم منعکس می نماید؟

اینها و بسیاری از سؤالات دیگر ما را با حرکتی رو به جلو مواجه می کند. در بخش های آینده این سند، ما نگاهی خواهیم داشت به بعضی از فنون فهم وب پویا که می تواند الهام بخش جامعه آرشیوسازی وب باشد. سپس به چالش های حرکت رو به جلویی که وب آرشیوی را به صورت بالقوه برای تحقیق، ارزشمندتر می کند، اشاره خواهیم نمود.

یادگیری از وب پویا

باید پرسید که چرا در دنیای پژوهشی اینترنت، آرشیوهای وب به صورت شهروندان درجه دوم به نظر می رسند؟ در مقایسه با کسانی که در وب پویا مطالعه می کنند، به مراتب تعداد کمتری از پژوهشگران از آرشیوهای وب استفاده می کنند و تعداد کمی از ابزارهای غیر تخصصی برای آرشیوهای وب در حال ساخت هستند، به ویژه در مقایسه با ابزارهایی که برای مطالعه وب پویا ساخته می شوند.

چالش عمومی برجسته در این بخش این است که، جامعه آرشیوسازی وب به ایجاد ارتباط بین منابعی که تولید می کنند با ابزارهای لبه برش (2) نیاز دارند، که توسط متخصصان رایانه، پژوهشگران، توسعه

نداریم، همچنان که این مورد فراتر از حیطه اختیار و تخصص ماست.

2- (نام نوعی فناوری) cutting edge

دهندگان مستقل، و هکرها برای مطالعه وب پویا گسترش یافته اند.

در حال حاضر، انواع ابزارهای توسعه یافته برای مطالعه وب پویا، روی هم رفته، به آسانی برای مطالعه داده های قابل دسترسی آرشیوهای وب به کار گرفته می شوند. این [امر] بزرگ ترین مانع برای فهم وب، نه تنها به عنوان یک تصویر لحظه ای بلکه به عنوان یک اکوسیستم در حال توسعه است.

مصورسازی

هر آرشیوی که ایجاد شده و مورد استفاده قرار می گیرد به [مرور] وسیع، غیر قابل نظر اجمالی و بدون نقشه یا دسترسی بصری و روش حسی برای مشاهده آرشیوها و چگونگی پیوندهای آن ها، خواهد شد. مصورسازی در اینجا یک راه حل اساسی است، اما انواع متعددی از آن ها در دسترس است.

چالش ها

ابزارهای بسیاری برای بررسی ارتباط بین کاربران رسانه های اجتماعی، ابزارهایی برای مرور پیوستار زمانی، ترکیب نقشه ها با ردپای کاربران، ابزارهای تصویری برای نمایش نحوه پیوند داده ها و نظیر آن وجود دارد؛ اما اینها به منظور کار کردن با آرشیوهای وب نیاز به متمرکز شدن دارند. مصورسازی اطلاعات حوزه پیشرفته پژوهشی است، اما بهترین نحوه مشاهده یک آرشیو (یا جست و جوی یکی از آن ها، یا دیدن ارتباطات درونی و ما بین آن ها) شامل رابط های حسی (بصری)، تغییرات دیداری تصاویر سه بُعدی و پوششگران زمانی، هنوز دست نیافتنی هستند؛ و برای مثال آیا احتمال شناسایی موضوع هایی درون مجموعه به وسیله روش های بازبینی تصویری یا سازماندهی تصویری یک سری از ارتباطات وجود دارد؟ به بیان دیگر، ایجاد ابزارهای مصورسازی برای کمک به پژوهشگران ممکن است؟

مثال ها: (2) Touch graph، (1) Apple Time Machine

عکس

دهندگان مستقل، و هکرها برای مطالعه وب پویا گسترش یافته‌اند.

در حال حاضر، انواع ابزارهای توسعه یافته برای مطالعه وب پویا، روی هم رفته، به آسانی برای مطالعه داده‌های قابل دسترسی آرشیوهای وب به کار گرفته می‌شوند. این [امر] بزرگ‌ترین مانع برای فهم وب، نه تنها به‌عنوان یک تصویر لحظه‌ای بلکه به‌عنوان یک اکوسیستم در حال توسعه است.

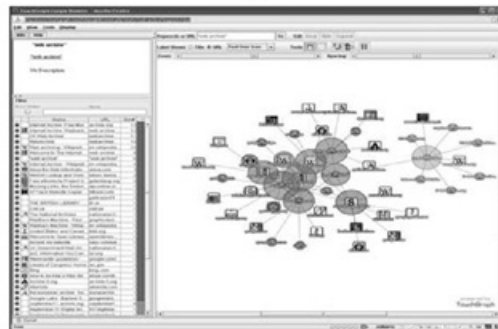
مصورسازی

هر آرشیوی که ایجاد شده و مورد استفاده قرار می‌گیرد به [مرور] وسیع، غیرقابل نظر اجمالی و بدون نقشه یا دسترسی بصری و روش حسی برای مشاهده آرشیوها و چگونگی پیوندهای آنها، خواهد شد. مصورسازی در اینجا یک راه‌حل اساسی است، اما انواع متعددی از آنها در دسترس است.

چالش‌ها

ابزارهای بسیاری برای بررسی ارتباط بین کاربران رسانه‌های اجتماعی، ابزارهایی برای مرور پیوستار زمانی، ترکیب نقشه‌ها با ردپای کاربران، ابزارهای تصویری برای نمایش نحوه پیوند داده‌ها و نظیر آن وجود دارد؛ اما اینها به‌منظور کار کردن با آرشیوهای وب نیاز به متمرکز شدن دارند. مصورسازی اطلاعات حوزه پیشرفته پژوهشی است، اما بهترین نحوه مشاهده یک آرشیو (یا جست‌وجوی یکی از آنها، یا دیدن ارتباطات درونی و ما بین آنها) شامل رابط‌های حسی (بصری)، تغییرات دیداری، تصاویر سه بُعدی و پوششگران زمانی، هنوز دست نیافتنی هستند؛ و برای مثال، آیا احتمال شناسایی موضوع‌هایی درون مجموعه به‌وسیله روش‌های بازبینی تصویری یا سازماندهی تصویری یک سری از ارتباطات وجود دارد؟ به بیان دیگر، ایجاد ابزارهای مصورسازی برای کمک به پژوهشگران ممکن است؟

مثال‌ها: 'Touch graph'، 'Apple Time Machine'



تصویر ۱. Touchgraph، در اینجا، توانایی بررسی پیوندهای میان وبگاه‌ها با استفاده از رابط گرافیکی را ارائه می‌کند. داده‌ها از وب پویا ترسیم شده‌اند.

1. Touchgraph(<http://www.touchgraph.com/>)
2. Apple Time Machine (<http://www.apple.com/macosx/what-is-macosx/time-machine.html>)

تصویر 1. Touchgraph، در اینجا، توانایی بررسی پیوندهای میان وب گاه‌ها با استفاده از رابط گرافیکی را ارائه می‌کند. داده‌ها از وب پویا ترسیم شده‌اند.

برنامه های کاربردی جست و جو همانند شکارچی

همان طوری که اطلاعات اینترنت به تکثیر و افزایش، هم در حجم و هم در تنوع انواع محتوا ادامه می دهند، جست و جوهای به مراتب پیچیده تری برای قابلیت استخراج هر چیز با معنا و استفاده از این مجموعه عظیم مورد نیاز خواهد بود. جست و جو مانند جست و جوی تصویر و ویدئو به سوی تکلیفی پیچیده تر سوق می یابد.

چالش: ایجاد سطح بلند پروازانه تری از برنامه ها با هزینه ای تقریباً ناچیز، به ویژه ابزارهایی که برای مجموعه ها می توانند به کار گرفته شوند. این، امر ممکن است نیازمند طراحانی برای استفاده جسورانه از موتورهای جست و جوی مبتنی بر فناوری ابر و زبان های جست و جوی بهتر با قابلیت انعطاف در پرسش سؤال های پیچیده از داده های موجود در آرشیوهای وب، باشد.

مثال ها: (در حال حاضر آپاچی) یا هوا! طرح زیربنایی [1](http://pig.apache.org/) (PIG Latin) برای پشتیبانی موردی تحلیل سری داده های خیلی بزرگ (تصویر 2)

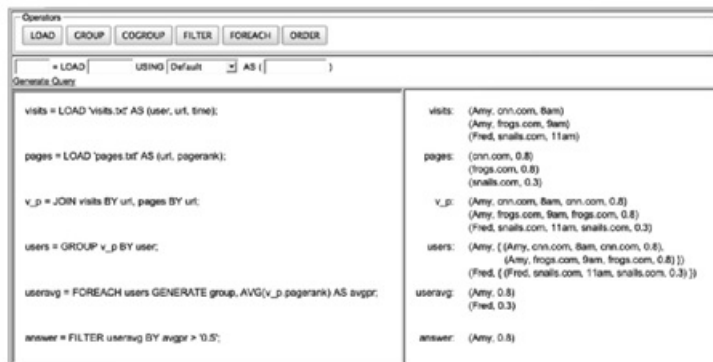
عکس

برنامه‌های کاربردی جست‌وجو همانند شکارچی

همان‌طوری که اطلاعات اینترنت به تکثیر و افزایش، هم در حجم و هم در تنوع انواع محتوا ادامه می‌دهند، جست‌وجوهای به مراتب پیچیده‌تری برای قابلیت استخراج هر چیز با معنا و استفاده از این مجموعه عظیم مورد نیاز خواهد بود. جست‌وجو مانند جست‌وجوی تصویر و ویدئو به سوی تکلیفی پیچیده‌تر سوق می‌یابد.

چالش: ایجاد سطح بلند پروازانه تری از برنامه‌ها با هزینه‌ای تقریباً ناچیز، به‌ویژه ابزارهایی که برای مجموعه‌ها می‌توانند به‌کار گرفته شوند. این امر، ممکن است نیازمند طراحی برای استفاده جسورانه از موتورهای جست‌وجوی مبتنی بر فناوری ابر و زبان‌های جست‌وجوی بهتر با قابلیت انعطاف در پرسش سؤال‌های پیچیده از داده‌های موجود در آرشیوهای وب، باشد.

مثال‌ها: (در حال حاضر آپاچی) یا هو! طرح زیربنایی^۱ PIG Latin برای پشتیبانی موردی تحلیل سری داده‌های خیلی بزرگ (تصویر ۲)



تصویر ۲. نمایی از Pi6 Pen نشان دهنده برنامه‌ای است که افراد متقاضی مشاهده صفحه‌های با رتبه‌بندی بالا را می‌یابد. (اولستون^۲، رید^۳، سریواستاوا^۴، کومر^۵ و تامکینز^۶، ۲۰۰۸)

تحلیل‌های شبکه اجتماعی

تحلیل شبکه اجتماعی (SNA) حوزه قابل توجهی از علاقه و فعالیت پژوهشی در میان پژوهشگران اینترنت، جامعه‌شناسان، فیزیکدانان و بسیاری دیگر است. تنوع موضوع‌های مورد توجه بسیار گسترده

1. PIG Latin(<http://pig.apache.org/>)
2. Olston
3. Reed
4. Srivastava
5. Kummer
6. Tomkins

تصویر ۲. نمایی از Pi6 Pen نشان دهنده برنامه‌ای است که افراد متقاضی مشاهده صفحه‌های با رتبه‌بندی بالا را می‌یابد. (اولستون^۲، رید^۳، سریواستاوا^۴، کومر^۵ و تامکینز^۶، ۲۰۰۸)

تحلیل‌های شبکه اجتماعی

تحلیل شبکه اجتماعی (SNA) حوزه قابل توجهی از علاقه و فعالیت پژوهشی در میان پژوهش‌گران اینترنت، جامعه‌شناسان، فیزیک‌دانان و بسیاری دیگر است. تنوع موضوع‌های مورد توجه بسیار گسترده

PIG Latin -1

Olston -2

Reed -3

Srivastava -4

Kummer -5

Tomkins -6

است که شامل فهم ارتباطات بین دوستان در شبکه های اجتماعی مانند فیس بوک (هوگان 2010 (1)) ، بررسی وابستگی های سیاسی مشارکت کنندگان در مباحث سیاسی (هیندمن (2)) ، (2007) ، و کشف شبکه های توطئه (طرح) در ادبیات انگلیسی (مورتنی 2011 ، (3) 2005) است.

ابزارهای جست و جوی مبتنی بر تحلیل های شبکه اجتماعی شامل (7) NodeXL (6) Voson ، (5) Pajec ، (4) UCINET ، و بسیاری دیگر، در حال تکثیر و توسعه هستند.

اگر چه تعداد کمی از این ابزارها، هر چه که باشد، برای استفاده به وسیله آرشيوهای وب توانمند یا بهینه شده اند.

چالش: نخست، کار با طراحان اصلی ابزارهای تحلیل شبکه اجتماعی برای توانمند و بهینه سازی آن ها در کار با داده های آرشيوی وب.

همچنین، توسعه روش های جدید ابداعی با احتمال یک بار در بعد زمان، که به داده های شبکه برای ردیابی مواردی همچون تکامل تدریجی شبکه های اجتماعی در طول زمان، افزوده می شود. به وسیله آرشيوسازی، نه تنها وضعیت سایت های شبکه های اجتماعی، بلکه زمانی که افراد پیوندها را ایجاد، نگهداری و حذف نموده، با دیگران ارتباط برقرار کرده، به گروه ها می پیوندند و گروه ها یا وبگاه ها را ترک می کنند.

ما باید به خاطر داشته باشیم که وب شبکه ای از پیوندهاست و تحلیل وب بینشی، کلی به ماهیت آن شبکه برای ما فراهم می کند.

مثال ها:

تحلیل های فیس بوک: ابزارهای بسیاری برای تحلیل تعاملات بین کاربران فیس بوک، جریان نفوذ و منحنی اجتماعی در دسترس هستند.

نمایش مثالی از فیس بوک (8):

عکس

است که شامل فهم ارتباطات بین دوستان در شبکه‌های اجتماعی مانند فیس بوک (هوغان^۱، ۲۰۱۰)، بررسی وابستگی‌های سیاسی مشارکت کنندگان در مباحث سیاسی (هیندمن^۲، ۲۰۰۷)، و کشف شبکه‌های توطئه (طرح) در ادبیات انگلیسی (مورتی^۳، ۲۰۱۱، ۲۰۰۵) است.

ابزارهای جست‌وجوی مبتنی بر تحلیل‌های شبکه اجتماعی شامل، Voson^۴، Pajec^۵، UCINET^۶، NodeXL^۷، و بسیاری دیگر، در حال تکثیر و توسعه هستند.

اگر چه تعداد کمی از این ابزارها، هر چه که باشد، برای استفاده به‌وسیله آرشیوهای وب توانمند یا بهینه شده‌اند.

چالش: نخست، کار با طراحان اصلی ابزارهای تحلیل شبکه اجتماعی برای توانمند و بهینه‌سازی آنها در کار با داده‌های آرشیوی وب.

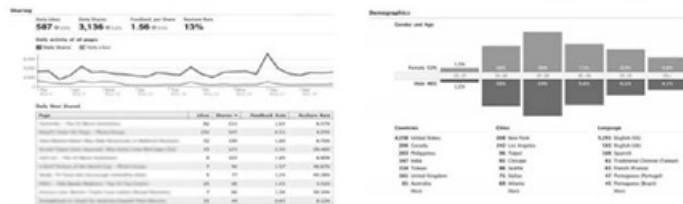
همچنین، توسعه روش‌های جدید ابداعی با احتمال یک‌بار در بعد زمان، که به داده‌های شبکه برای ردیابی مواردی همچون تکامل تدریجی شبکه‌های اجتماعی در طول زمان، افزوده می‌شود، به وسیله آرشیوسازی، نه تنها وضعیت سایت‌های شبکه‌های اجتماعی، بلکه زمانی که افراد پیوندها را ایجاد، نگهداری و حذف نموده، با دیگران ارتباط برقرار کرده، به گروه‌ها می‌پیوندند و گروه‌ها یا وبگاه‌ها را ترک می‌کنند.

ما باید به خاطر داشته باشیم که وب شبکه‌ای از پیوندهاست و تحلیل وب بینشی، کلی به ماهیت آن شبکه برای ما فراهم می‌کند.

مثالها:

تحلیل‌های فیس بوک: ابزارهای بسیاری برای تحلیل تعاملات بین کاربران فیس بوک، جریان نفوذ و منحنی اجتماعی در دسترس هستند.

نمایش مثالی از فیس بوک^۸:



تصاویر نمودارهای اجتماعی فیس بوک:

1. Hogan
2. Hindman
3. Moretti
4. <http://www.analytictech.com/ucinet/>
5. <http://pajec.imfm.si/doku.php>
6. <http://voson.anu.edu.au/>
7. <http://nodexl.codeplex.com/>
8. <http://www.facebook.com/insights/>

تصاویر نمودارهای اجتماعی فیس بوک:

[/http://www.analytictech.com/ucinet](http://www.analytictech.com/ucinet) -4

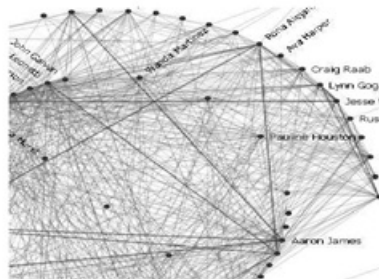
<http://pajek.imfm.si/doku.php> -5

[/http://voson.anu.edu.au](http://voson.anu.edu.au) -6

[/http://nodexl.codeplex.com](http://nodexl.codeplex.com) -7

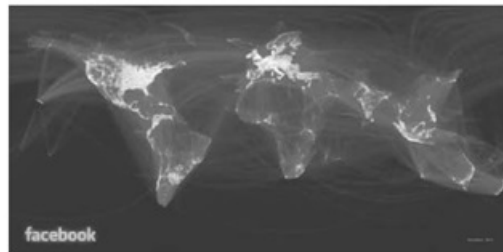
<http://www.facebook.com/insights/> -8

۱۲۰ مدیریت منابع اطلاعاتی وب



تصویر ۳. منبع:

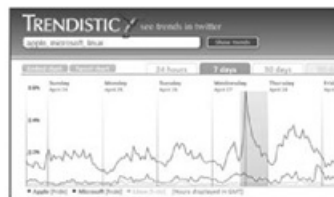
http://infosthetics.com/archives/2008/03/facebook_social_network_graph.html



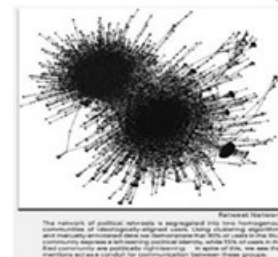
تصویر ۴. منبع:

<http://www.facebook.com/notes/facebook-engineering/visualizing-friendships/469716398919>

به همین ترتیب، طیف گسترده‌ای از تحلیل‌های توییت‌ر مانند Twitalyzer^۱، Trendistic^۲، و دیگران وجود دارد.



تصویر ۵. trndistic.



تصویر ۶. Truthy. (<http://truthy.indiana.edu/>)

1. <http://www.twitalyzer.com/>
2. <http://trendistic.com/>

تصویر ۳. منبع: http://infosthetics.com/archives/2008/03/facebook_social_network_graph.html

تصویر ۴. منبع: <http://www.facebook.com/notes/facebook-engineering/visualizing-friendships/469716398919>

به همین ترتیب، طیف گسترده ای از تحلیل های توییتز مانند [\(2\) Twitalyzer \(1\)](#)، [Trendistic](#)، و دیگران وجود دارد.

تصویر 5. trndistic.

تصویر 6. Truthy (<http://truthy.indiana.edu/>)

ص: 120

1- [/http://www.twitalyzer.com](http://www.twitalyzer.com)

2- [/http://trendistic.com](http://trendistic.com)

Far Left	Moderate Left	Center	Moderate Right	Far Right
#healthcare	#arp #women	#democrats #social	#rangel #vaste	#912project #twisters
#judaism #hollywood	#citizensunited	#seniors #dnc	#saveamerica	#gop2112 #israel
#2010elections	#democratic	#budget #political	#american #gold	#foxnews #mediabias
#capitalism #recession	#banksters #energy	#goproud #christian	#repeal #mexico	#constitution
#security #dceasact	#sarahpalin	#media #nobel	#terrorian #gopleader	#patriots #cednov
#publicoption	#progressives		#palin12	#abortion
#topprogs	#stopbeck #iraq			

تصویر ۸. برچسب‌های نشانه‌گذاری شده به وسیله کاربران حوزه سیاست (کونور^۱، راتکوویچ^۲، فرانسیسکو^۳، و دیگران: ۲۰۱۱)



تصویر ۷. توییت کردن و سیاست (کونور، راتکوویچ، گونکالوز، فلامینی و منسز^۴، ۲۰۱۱)

سنجش‌های دگرساز

اصطلاحی در حال پدیدار شدن برای راه‌های جدید اندازه‌گیری تأثیر علمی، فراتر از معیارهای کتابسنجی، وب‌سنجی و علم‌سنجی است.

ارتباطات بین ما و بین دانشمندان و درون گروه‌های علمی به‌طور روزافزونی بر روی وب در حال شکل‌گیری است. جامعه در حال ظهوری از پژوهشگران که مطالعه پژوهشی انجام می‌دهند، از توالی‌ها و پیوندهایی که به وسیله ابزارهایی چون توییت، مندلی^۵، فرند فید و بسیاری دیگر بر جای گذاشته شده است، برای درکی که بسیار سریع‌تر از تأثیرات سنتی می‌توانند توسعه یابند، استفاده می‌شوند.

فراتر از چیزی که می‌توان به‌عنوان سنجش دگرساز تصور کرد: چگونگی ردیابی مشارکت‌های غیرآکادمیک به سمت دانش است.

آیا می‌توان ابزارهای تحقیقاتی پژوهشگران را برای سایر حوزه‌های غیر تخصصی به‌کار برد؟ برای مثال، آیا می‌توانیم تأثیر مشارکت کنندگان انفرادی بر گروه‌های سرگرمی در طول زمان را با استفاده از معیارهای مشابه برای درک چگونگی تحول تأثیر پژوهشگر، ارزیابی کنیم؟

چالش: فعال‌سازی روش‌هایی به مراتب ساده‌تر برای مشخص کردن طیف زمانی مواد رقومی، بنابراین، تحلیل سنجش دگرساز می‌تواند به روش‌های کاملاً مشابه کتابسنجی انجام شود: در الگوی

1. Conover
2. Ratkiewicz
3. Francisco
4. Goncalves
5. Flammim
6. Menczer
7. Mendeley

تصویر ۸. برچسب‌های نشانه‌گذاری شده به وسیله کاربران حوزه سیاست (کونور^(۱)، راتکوویچ^(۲)، فرانسیسکو^(۳)، و دیگران؛

(2011)

اصطلاحی در حال پدیدار شدن برای راه های جدید اندازه گیری تأثیر علمی، فراتر از معیارهای کتاب سنجی، وب سنجی و علم سنجی است.

ارتباطات بین ما و بین دانشمندان و درون گروه های علمی به طور روزافزونی بر روی وب در حال شکل گیری است. جامعه در حال ظهوری از پژوهشگران که مطالعه پژوهشی انجام می دهند، از توالی ها و پیوندهایی که به وسیله ابزارهایی چون توئیتر مندلی (4)، فرند فید و بسیاری دیگر بر جای گذاشته شده است، برای درکی که بسیار سریع تر از تأثیرات سنتی می توانند توسعه یابند، استفاده می شوند.

فراتر از چیزی که می توان به عنوان سنجش دگر ساز تصور کرد: چگونگی ردیابی مشارکت های غیرآکادمیک به سمت دانش است.

آیا می توان ابزارهای تحقیقاتی پژوهشگران را برای سایر حوزه های غیر تخصصی به کار برد؟ برای مثال، آیا می توانیم تأثیر مشارکت کنندگان انفرادی بر گروه های سرگرمی در طول زمان را با استفاده از معیارهای مشابه برای درک چگونگی تحول تأثیر پژوهش گر، ارزیابی کنیم؟

چالش: فعال سازی روش هایی به مراتب ساده تر برای مشخص کردن طیف زمانی مواد رقومی، بنابراین، تحلیل سنجش دگر ساز می تواند به روش های کاملاً مشابه کتاب سنجی انجام شود: در الگوی

ص: 121

Conover -1

Ratkiewicz -2

Francisco -3

Mendeley -4

انتشارات رسمی، هر نشریه ای یک نویسنده و تاریخ نشر دارد و برای ردیابی توالی استنادها به یک اثر انفرادی استفاده می شود. نشریات، نظام آرشیوی دارند که آزمایش شده و به نحو قابل اعتماد و مناسبی پذیرفته شده اند: نشریه تخصص با وجود این، نشریات غیر رسمی وب، روش توسعه یافته مناسب و مشابهی برای آرشیو سازی سهم ها (مشارکت ها) با دانش، به شکلی که قابل استناد باشند و در طی زمان دوباره جایگزین شوند، ندارند. این شکافی است که انتظار می رود رفع شود و آرشیوهای وب زمینه های بارزی برای شروع هستند.

همچنین، برای مشارکت های غیر تخصصی چه ابزارهایی برای تحلیل ویکی ها و موجودیت های مشترک استفاده می شود که می تواند فهم این تغییرها را در دوره های زمانی توسعه دهد؟

مثال ها: [\(1\) Data Cite](#) و [\(2\) Reader Meter](#)

عکس

انتشارات رسمی، هر نشریه‌ای یک نویسنده و تاریخ نشر دارد و برای ردیابی توالی استنادها به یک اثر انفرادی استفاده می‌شود. نشریات، نظام آرشیوی دارند که آزمایش شده و به نحو قابل اعتماد و مناسبی پذیرفته شده‌اند: نشریه تخصصی با وجود این، نشریات غیر رسمی وب، روش توسعه یافته مناسب و مشابهی برای آرشیو سازی سهم‌ها (مشارکت‌ها) با دانش، به شکلی که قابل استناد باشند و در طی زمان دوباره جایگزین شوند، ندارند. این شکافی است که انتظار می‌رود رفع شود و آرشیوهای وب زمینه‌های بارزی برای شروع هستند.

همچنین، برای مشارکت‌های غیر تخصصی، چه ابزارهایی برای تحلیل ویکی‌ها و موجودیت‌های مشترک استفاده می‌شود که می‌تواند فهم این تغییرها را در دوره‌های زمانی توسعه دهد؟

مثال‌ها: ¹Reader Meter و ²Data Cite



تصویر ۲. reader meter نوعی از ابزار Alt metric برای درک نحوه خواندن یک نویسنده، بر اساس

آمارهای مندی (<http://www.mendeley.com/>) منبع: <http://readermeter.org>

وب‌نوشت (حاشیه نگاری) اجتماعی

کاربران مایل اند که بتوانند پیوندها و برچسب‌هایشان^۳ و همچنین نظرها و حاشیه‌نویسی‌هایشان را بر روی منابع نشر دهند. برای پژوهشگران، فهم چگونگی توسعه این جوامع در طول زمان و نگهداری خودشان سؤال مهمی است. برای مثال Reddit بیش از ۸ میلیون خواننده انحصاری و یک میلیون بازدید صفحه در هر ماه دارد (جسرا، ۲۰۱۱).

چالش‌ها: گستره‌ای راکه آرشیوها می‌توانند نه فقط وبگاه‌ها و مجموعه‌هایشان، بلکه پیوندها و

1. ReaderMeter
2. Data Cite
3. bookmark

تصویر ۲. reader meter نوعی از ابزار Alt metric برای درک نحوه خواندن یک نویسنده بر اساس آمارهای مندی (<http://www.mendeley.com/>) منبع: <http://readermeter.org>

وب‌نوشت (حاشیه نگاری) اجتماعی

کاربران مایل اند که بتوانند پیوندها و برچسب‌هایشان (3) و همچنین نظرها و حاشیه‌نویسی‌هایشان را بر روی منابع نشر دهند. برای پژوهشگران فهم چگونگی توسعه این جوامع در طول زمان و نگهداری خودشان سؤال مهمی است. برای مثال Reddit بیش از ۸ میلیون

خواننده انحصاری و یک میلیون بازدید صفحه در هر ماه دارد (جسرا، 2011).

چالش‌ها: گستره ای را که آرشیوها می توانند نه فقط وب گاه ها و مجموعه های شان، بلکه پیوند ها و

ص: 122

Data Cite -1

ReaderMeter -2

bookmark -3

حاشیه نویسی های آن صفحه ها و مجموعه ها را ذخیره نمایند، در نظر بگیرید. توانایی پاسخ به این سؤال که «چگونه افراد این منابع را برای یکدیگر نشانه گذاری می نمایند، و چگونه در طول زمان تغییر می کند؟» و بررسی استفاده از فناوری های موجود ترکیب و تطبیق دادن ابزارهای حاشیه نویسی اجتماعی موجود. گامی بیشتر بردارید، آیا می توان اجتماعی از پیوندها و وب نوشت ها برای مجموعه های آرشیوی با استفاده از ابزارهای مشابه مورد تأکید قرار داد تا معلوم شود افراد چگونه موارد موجود روی وب پویا را برای یکدیگر نشانه گذاری می کنند؟

مثال ها: (1) Delicious، ردایت بر مبنای بوک مارکلت (2) مثال: (3) Madcow)

معماران جدید

برای استفاده از داده های وب به منظور درک اتکای دنیا بر منابعی چون (4) API، و داده های پیوند شده، فعالیتی تعریف شده است تا داده ها را برای استفاده و هدف گذاری مجدد و ترکیب [آماده سازی]، قابل نماید. دلیل مهمی که چرا وب پویا نسبت به وب آرشیوی بسیار فعالانه تر مورد جست و جو قرار می گیرد، این است که طراحان ابزار به داده های وب پویا یا از طریق خزش گره های وبگاه ها به طور مستقیم و یا از طریق API ها به گوگل / یاهو، توئیتر، فیسبوک و مانند آن می توانند دسترسی یابند. حتی با وجود محدودیت بعضی از این API ها دلیل اصلی، پژوهش هایی هستند که با استفاده از وب پویا رونق گرفته اند.

API روش قدرتمندی برای ساخت نرم افزار های کاربردی جدید است که بر روی داده ها طراحی شده و منابع چندگانه داده را ترکیب می نماید.

پژوهشگری به ما گفت: من در مورد استفاده از برچسب نشانه گذاری اصلی در توئیتر تحقیق می کنم، و محدودیت شان را در استفاده از API در حالتی که بیشترین آشفستگی را در حین توئیت کردن دارند می یابم، من می خواهم به آن ها در حالی که هنوز برخط و در دسترس هستند دسترسی داشته باشم، اگر چه استقرار آن ها، یا به عبارت دیگر اجرای گزارش ها یا انباشت آن ها کاملاً مشکل است. برای مثال آن ها Twapper keeper را محدود به خدمت قابل دسترسی کردند که اجازه آماده سازی گزارش هایی برای کار در مورد این برچسب ها که از آن ها شناخت داشتیم را می داد.

سؤال بعدی این است که چگونه API ها می توانند آرشیو شوند و چه زمانی محدود یا متوقف می شوند؟ چگونه منابعی که تاکنون آرشیو شده اند از API تأثیر می گیرند؟ آیا روش هایی برای حفاظت محتوا همراه با حفظ و احترام به حقوق مالکان و دوره های مجوز وجود دارد؟

چالش: چگونه داده های وب قدیمی از طریق API ها باز شده و پیوند داده می شوند، تا افراد هوشمند خارج از آن بتوانند برای کاربرد راه های جدید تولید دانش آن ها را ترکیب کنند [؛ البته] به وسیله ارائه ابزارهایی به افراد، برای تعریف روال های کاری به صورتی انعطاف پذیرتر از اینکه وادار به استفاده از

ص: 123

bookmarklet -2

(/MadCow(<http://www.web-notes.com> -3

-4) (رابط یا میانجی برنامه های کاربردی) API: Application Program Interface

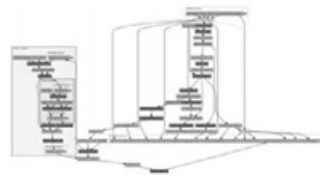
ابزارهایی یک منظوره یکپارچه بشوند. یک رهیافت، پیاده سازی کارکردهای تحلیلی مانند خدمات وب خواهد بود که می تواند با گردش کار قانونی موتور ترکیب شود. این رهیافت به صورت گسترده ای در bioinformatic که با مشکلات مشابهی در خصوص یکپارچه سازی داده ها از مخازن متعدد رویه روست استفاده می شود.

مثال: گردش کار قانون - مصوب موتور [1 Taverna](#) برای ترکیب [تلفیق] خدمات وب پیشنهادی توسط پردازش توزیع شده و نظام های ذخیره سازی استفاده می شود، سپس گردش کارها می توانند به اشتراک در آمده، دوباره استفاده و هدفگذاری شوند. myexperience [2](#) مثالی از یک مخزن قابل اشتراک گردش کارهای علمی:

عکس

ابزارهایی یک منظوره یکپارچه بشوند. یک رهیافت، پیاده‌سازی کارکردهای تحلیلی مانند خدمات وب خواهد بود که می‌تواند با گردش کار قانونی موتور ترکیب شود. این رهیافت به صورت گسترده‌ای در bioinformatic که با مشکلات مشابهی در خصوص یکپارچه سازی داده‌ها از مخازن متعدد روبه‌روست، استفاده می‌شود.

مثال: گردش کار قانون - مصوب موتور Taverna^۱ برای ترکیب [تلفیق] خدمات وب پیشنهادی توسط پردازش توزیع شده و نظام‌های ذخیره‌سازی استفاده می‌شود، سپس گردش کارها می‌توانند به اشتراک در آمده، دوباره استفاده و هدفگذاری شوند. ^۲ myexperience مثالی از یک مخزن قابل اشتراک گردش کارهای علمی:



تصویر ۱۱. گردش کارهای تاورنا



تصویر ۱۰. منبع: کارول گلوبال و دیگران^۲

ماشین‌های اجتماعی

تیم برنزی و همکارانش در مورد اینکه وب در حال تبدیل شدن به یک ماشین اجتماعی است، بحث کرده‌اند. به این معنی که فقط یک مخزن اطلاعات نیست، بلکه زیرساختی برای همکاری در رفع مشکل اجرای وظایف موجودیت‌های انسانی است که به آسانی نمی‌توانند توسط ماشین انجام شوند. دانشمندان علوم اجتماعی علاقه‌مند به درک تعامل بین فناوری و علوم اجتماعی، می‌خواهند بدانند که چگونه افراد و فناوری برای حل وظایف پیچیده‌ای که هر یک به تنهایی از عهده رفع آن بر نمی‌آیند، با یکدیگر همکاری می‌کنند.

چالش: چگونه می‌توانیم تجربه و تعامل بین کاربران و ماشین‌های اجتماعی بر روی وب را ذخیره و درک کنیم؟ همه امور اجتماعی در مورد تعامل‌هاست. اگر نتوانیم تعاملات را درک کنیم، هرگز نمی‌توانیم آنچه را که درباره ماشین اجتماعی است درک کنیم.

مثال: ^۳ Amazon's Mechanical Turk مکانیسمی برای توزیع مشکلات بین انسان‌های خیره و جمع‌آوری راه‌حل‌هاست. گردآوری منابع، می‌تواند برای رفع مشکلات سخت به کار رود. به‌عنوان مثال

1. Taverna (<http://www.taverna.org.uk/>)

2. <http://www.myexperiment.org/>

3. Carol Global (http://nar.oxfordjournals.org/content/38/suppl_2/W677.full)

4. Amazon's Mechanical Turk (<https://www.mturk.com/mturk/welcome>)

تصویر 11. گردش کارهای تاورنا

تصویر 10: منبع: کارول گلوبال و دیگران (3)

ماشین‌های اجتماعی

تیم برنزی و همکارانش در مورد اینکه وب در حال تبدیل شدن به یک ماشین اجتماعی است، بحث کرده‌اند. به این معنی که فقط یک

مخزن اطلاعات نیست، بلکه زیرساختی برای همکاری در رفع مشکل اجرای وظایف موجودیت های انسانی است که به آسانی نمی توانند توسط ماشین انجام شوند. دانشمندان علوم اجتماعی علاقه مند به درک تعامل بین فناوری و علوم اجتماعی، می خواهند بدانند که چگونه افراد و فناوری برای حل وظایف پیچیده ای که هر یک به تنهایی از عهده رفع آن بر نمی آیند، با یکدیگر همکاری می کنند.

چالش: چگونه می توانیم تجربه و تعامل بین کاربران و ماشین های اجتماعی بر روی وب را ذخیره و درک کنیم؟ همه امور اجتماعی در مورد تعامل هاست. اگر نتوانیم تعاملات را درک کنیم، هرگز نمی توانیم آن چه را که درباره ماشین اجتماعی است درک کنیم.

مثال: [Amazon's Mechanical Turk](#) (4) مکانیسمی برای توزیع مشکلات بین انسان های خبره و جمع آوری راه حل هاست گردآوری منابع می تواند برای رفع مشکلات سخت به کار رود. به عنوان مثال

ص: 124

1 - <http://www.taverna.org.uk> (/Taverna)

2 - <http://www.myexperiment.org>

3 - http://nar.oxfordjournals.org/content/38/suppl_2/W677.full (Carol Global)

4 - <https://www.mturk.com/mturk/welcome> (Amazon's Mechanical Turk)

1) CAPTCHA برای شناسایی نویسه نوری کمک کننده انسان است. افراد چگونه با این ابزارها و از طریق اینها با یکدیگر تعامل می کنند؟

برای مثال، در مورد «بازهای هدف دار» همچون (2) Foldit، چالش نه در آرشیوسازی وبگاه ها و نه فقط در بازی بلکه در چگونگی بازی افراد است. کاربران چگونه با بازی تعامل دارند؟ از تحقیق در مورد چگونگی بازی بازیکنان در تعامل با فضای پیوسته، دروسی را می توان طراحی نمود (برای مثال: ویلیامز (3)، یی (4) و کاپلان (5)، 2008).

شبکه های نقشه برداری

به طور روزافزونی جغرافی دانان از داده های اینترنت برای دریافت اطلاعات مکان ها، جریان ها، جهات و ثروت، فقر و تغییر شکل محتوا، و تأثیر در طول زمان و فضا استفاده می کنند.

چالش: استخراج خودکار اطلاعات جغرافیایی از پیوندهای درونی و بیرونی درون یک مجموعه که می تواند نقشه برداری شود. این چالش، در حال حاضر در وب پویا وجود دارد و حتی در زمان افزایش پیچیدگی تغییرات در طول زمان، بیشتر هم می شود. در حال حاضر، بخش عمده ای از اطلاعات که می تواند با استفاده از روش های دو بعدی نمایان شود به داده های سه بعدی و چهار بعدی (همانند اسلایدر (6) ها) نیاز خواهند داشت. تا اطلاعات جغرافیایی را که در طول زمان تغییر می یابند، ایجاد کند برای مثال، درک تأثیر جغرافیایی درون و ما بین دانشگاه ها، دولت ها، و شرکت ها در طول زمان از نظر تئوری محتمل است، اما به استخراج اطلاعات جغرافیایی از داده های ساختار نیافته وب نیاز دارد.

مثال: (7) Floating Sheep

ص: 125

1 - CAPTCHA: برنامه ای کامپیوتری است برای اجرای آزمون های نرم افزاری که فقط انسان می تواند به آن پاسخ دهد. (دسترسی در <http://www.google.com/recaptcha/captcha>)

2 - Foldit

3 - Williams

4 - Yee

5 - Caplan

6 - Slider: رابط های کاربر گرافیکی برای ارتباط افراد با نرم افزار در ابزارهایی نظیر موبایل (دسترسی در <http://www.wisegeek.com/in-computing-what-is-a-slider.htm>)

7 - Floating Sheep (<http://www.floatingsheep.org>)



تصویر ۱۳. نقشه فلوتینگ شیپ «نقشه جغرافیایی مذهبی گوگل»

علم وب

علم وب تلاشی است که توسط محققان برای مطالعه وب به عنوان یک «مصنوع اطلاعاتی»، فهم چگونگی رشد و تکامل آن و کیفیت توسعه جوامع آن، انجام می‌گیرد.

چالش‌ها: نیاز به ابزارهای قدرتمندی برای تحلیل «نمودار وب» به عنوان شیء ریاضی. مکان شناسی آن چیست؟ چگونه گروه‌ها تکامل می‌یابند؟ چه نوع رتبه‌بندی قانونی به کار می‌رود؟ آیا واقعاً وب به وسیله قانونی قوی مدیریت می‌شود؟ چگونه اطلاعات روی وب منتشر می‌شوند؟

وب تنها یک فضای اطلاعاتی نیست، بلکه مجموعه‌ای پیچیده و خانواده‌ای از روابط درونی فضاهای فرعی است که محتوای اطلاعاتی آن در بعضی مواقع توسط جوامع مجزا تعیین می‌شوند. چگونه اطلاعات، بین این صفحه‌ها به اشتراک گذاشته و منتشر می‌شوند؟

برای پاسخ به این سؤال، به توسعه ابزارهایی نیاز داریم که قادر به ردیابی تکامل و انتقال مفاهیم در طول زمان و فضاهای مختلف باشند (برای مثال: بین وبلاگ‌نویسی و رسانه «mainstream»).

مثال: ^۱ Media Cloud و ^۲ Recorded future

1. Media Cloud (<http://cyber.law.harvard.edu/research/mediacloud>)

2. Recorded future (<https://www.recordedfuture.com/>)

تصویر 13. نقشه فلوتینگ شیپ «نقشه جغرافیایی مذهبی گوگل»

علم وب تلاشی است که توسط محققان برای مطالعه وب به عنوان یک «مصنوع اطلاعاتی»، فهم چگونگی رشد و تکامل آن و کیفیت توسعه جوامع آن، انجام می گیرد.

چالش ها: نیاز به ابزارهای قدرتمندی برای تحلیل «نمودار وب» به عنوان شیء ریاضی مکان شناسی آن چیست؟ چگونه گروه ها تکامل می یابند؟ چه نوع رتبه بندی قانونی به کار می رود؟ آیا واقعاً وب به وسیله قانونی قوی مدیریت می شود؟ چگونه اطلاعات روی وب منتشر می شوند؟

وب تنها یک فضای اطلاعاتی نیست، بلکه مجموعه ای پیچیده و خانواده ی از روابط درونی فضاهای فرعی است که محتوای اطلاعاتی آن در بعضی مواقع توسط جوامع مجزا تعیین می شوند. چگونه اطلاعات، بین این صفحه ها به اشتراک گذاشته و منتشر می شوند؟

برای پاسخ به این سؤال، به توسعه ابزارهایی نیاز داریم که قادر به ردیابی تکامل و انتقال مفاهیم در طول زمان و فضاهای مختلف باشند (برای مثال: بین ویلاگ نویسی و رسانه «mainstream»).

مثال: (1) Media Cloud و (2) Recorded future

ص: 126

1- (Media Cloud (<http://cyber.law.harvard.edu/research/mediacloud>

2- (/Recorded future (<https://www.recordedfuture.com>

آینده آرشیو وب ۱۲۷



تصویر ۱۴. Recorded Futur سیر اخبار را در طی زمان را دنبال می کند.

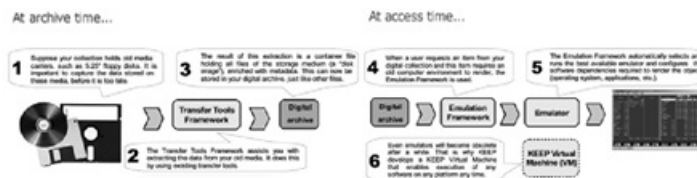
تصویر ۱۵. Media Cloud سیر اخبار جغرافیایی را دنبال میکند.

درک تجربه به جای محتوا

به طور فزاینده، دانشمندان به سمت درک اهمیت چگونگی استفاده افراد از محتوای وب و نه فقط محتوا به خودی خود، سوق پیدا کرده اند. این امر، به حساب وضعیت تجربه وب و محتوای قابل اجرا گذاشته می شود: تجربه درباره اینکه کدام طرح، کدام مرورگر یا افزونه ها یا رمز گذارهای دوطرفه مورد استفاده قرار می گیرد و به طور فزاینده به موقعیت مکانی کاربران وابسته است.

چالش ها: برای درک تجربه ها به توانایی تکرار تجربه نیازمندیم، طرح های زیربنایی، سیستم های عملیاتی، مرورگرها، و به همین ترتیب تغییر تجربه وب.

مثالها: Browsershots^۲ و KEEP^۳ (نگهداری نمونه سازی محیط های قابل انتقال)



تصویر ۱۶. KEEP (نگهداری نمونه سازی محیط های قابل انتقال)

تحلیل وب معنایی و مجموعه داده های پیوند شده

مجموعه داده های پیوند شده به سرعت با حداقل ۲/۷۵ بیلیون پیوند سه تایی در مجموعه های شناخته شده در

1. plugins
2. Browsershots (<http://browsershots.org>)
3. KEEP (<http://www.keep-project.eu>)

تصویر ۱۴. Recorded Futur سیر اخبار را در طی زمان را دنبال می کند.

تصویر ۱۵. Media Cloud سیر اخبار جغرافیایی را دنبال می کند.

به طور فزاینده دانشمندان به سمت درک اهمیت چگونگی استفاده افراد از محتوای وب و نه فقط محتوا به خودی خود، سوق پیدا کرده اند. این امر، به حساب وضعیت تجربه وب و محتوای قابل اجرا گذاشته می شود: تجربه درباره اینکه کدام طرح کدام مرورگر یا افزونه ها (1) یا رمزگذارهای دوطرفه مورد استفاده قرار می گیرد و به طور فزاینده به موقعیت مکانی کاربران وابسته است.

چالش ها: برای درک تجربه ها به توانایی تکرار تجربه نیازمندیم، طرح های زیربنایی، سیستم های عملیاتی، مرورگرها، و به همین ترتیب تغییر تجربه وب.

مثال ها: (2) Browsershots و (3) KEEP (نگهداری نمونه سازی محیط های قابل انتقال)

تصویر 16. KEEP (نگهداری نمونه سازی محیط های قابل انتقال)

تحلیل وب معنایی و مجموعه داده های پیوند شده

مجموعه داده های پیوند شده به سرعت با حداقل 28/5 بیلیون پیوند سه تایی در مجموعه های شناخته شده در

ص: 127

plugins -1

(Browsershots (<http://browsershots.org> -2

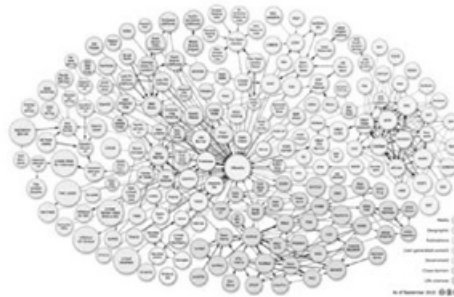
(KEEP (<http://www.keep-project.eu> -3

حال رشد هستند. ابزارها توسط وب معنایی یا جامعه داده های پیوند شده توسعه می یابند، به طوری که می توانند به سرعت مشکل مدیریت فراداده را آسان کرده و به این فراداده ها اجازه دهند که در مقیاس بزرگ قابل جست و جو بوده و همچنین برای یکپارچه سازی داده ها در مجموعه به روشی بسیار پیچیده تر استفاده شوند.

عکس

۱۲۸ مدیریت منابع اطلاعاتی وب

حال رشد هستند. ابزارها توسط وب معنایی یا جامعه داده های پیوند شده توسعه می یابند، به طوری که می توانند به سرعت مشکل مدیریت فراداده را آسان کرده و به این فراداده ها اجازه دهند که در مقیاس بزرگ قابل جست و جو بوده و همچنین برای یکپارچه سازی داده ها در مجموعه به روشی بسیار پیچیده تر استفاده شوند.



تصویر ۱۷. منبع: نمودار ابر گونه داده های باز پیوند خورده

توسط ریچارد کی گانیاک^۱ و انجا ینتس^۲ <http://lod-cloud.net/>

مثالها: ^۳ [Sig.ma](http://sig.ma/) / ^۴ جست و جوی مبتنی بر آر دی اف



تصویر ۱۸. [Sig.ma](http://sig.ma/) / [Sindic](http://sig.ma/) جست و جوی مبتنی بر RDF.

منبع: <http://sig.ma/search?q=Tim/۲۰Berners/۲۰Lee>

1. Richard Cyganiak
2. Anja Jentzsch
3. <http://sig.ma/>
4. Resource Description Framework (RDF)

تصویر 17. منبع: نمودار ابر گونه داده های باز پیوند خورده توسط ریچارد کی گانیاک (1) و انجا ینتس (2) http://lod-cloud.net

مثال‌ها: (3) [Sig.ma / Sindic](http://sig.ma) جست و جوی مبتنی بر آر دی اف (4)

تصویر 18. [Sig.ma / Sindic](http://sig.ma) جست و جوی مبتنی بر RDF

منبع: <http://sig.ma/search?q=Tim%20Berners%20Lee>

ص: 128

Richard Cyganiak -1

Anja Jentzsch -2

/http://sig.ma -3

Resource Description Framework (RDF) -4

اما گام‌های پیش رو برای جامعه آرشیوی وب چیست؟

بسیاری از مواردی که محققان [به آن‌ها] نیاز دارند، به وضوح قابل مشاهده هستند، ولی این به معنی در دسترس بودن آن‌ها نیست. در گزارش‌های قبلی GISC که با همکاران دیگر نوشتیم و در قسمت‌های پیشین توضیح دادیم (دوگرتی و دیگران 2010؛ توماس و دیگران، 2010) توجه‌های متعددی پیرامون موضوع اصلی ارائه دادیم:

ساخت جامعه، ساخت ابزار، و منابع و طراحی تجربه‌ها (دوگرتی و دیگران، 2010، ص 27-29). قصد نداریم که فهرستی کامل از پیشنهادها، گزارش‌های قبلی را در اینجا مطرح کنیم. 22 پیشنهاد در گزارش دوگرتی و همکارانش و بیش از 20 پیشنهاد در گزارش توماس و همکارانش وجود دارد. بنابراین، خواننده را به مرور این گزارش‌ها همانند این گزارش توصیه می‌کنیم با این حال می‌توانیم برای نیل به هدفمان بعضی از موضوع‌های مهمی را که محققان به آن‌ها علاقمند هستند، برجسته کرده و چالش‌هایی را شناسایی کنیم که آرشیوها برای کمک به آرشیوهای وب در تبدیل آن‌ها به بخشی از ابزارهای استاندارد برای محققان رشته‌های متنوع، با آن‌ها مواجه هستند.

این بخش از گزارش، موضوع‌ها و سؤال‌هایی را بر می‌شمارد که گروه‌های مختلف محققان می‌خواهند پرسند یا از آرشیو وب - آرشیو - در - جعبه - می‌خواهند پرسند. ما آن دسته را که چالش‌ها و راه‌حل‌های بالقوه‌ای دارند شناسایی کرده ایم بعضی از این راه‌حل‌ها، به طور ویژه راه‌حل‌هایی کوتاه مدت هستند که می‌توانند در سطح مؤسسه‌ها انجام شوند. بسیاری از رهیافت‌های نه‌چندان گسترده به نگاه وسیع‌تری در سطح ملی منطقه‌ای یا بین‌المللی توسط سازمان‌هایی چون IIPC نیاز دارند.

ما از آرشیو - در - جعبه یاد کردیم، زیرا در میان بعضی افراد این تفکر وجود دارد که وب آرشیوی فراتر از روزهای اولیه خودش در ایجاد صفحه‌های قابل دسترسی برای تحلیل‌های آینده، (تفکر تعامل سنتی (1) برای پایگاه Wayback Machin که به کاربر اجازه می‌داد به طور عمده به صفحه‌های منفرد آرشیو دسترسی داشته و آن‌ها را ببیند) در حال حرکت به سمت ایجاد مجموعه‌های قابل دسترسی به عنوان ابزارهای تحقیقی، است. برای مثال، زمانی که با دامنه «the UK government. Gov.uk» از 2011 - 2020 مواجه هستیم، یک پژوهشگر چه تصویری می‌تواند از توانایی انجام آن داشته باشد؟ چه سؤال‌هایی می‌تواند از یک مجموعه پرسیده شود؟ برای نمونه، محتوای کامل وب در مورد بانکدارهای اصلی وال استریت و سایت‌هایی که به طور مستقیم با آن پیوند دارند، در صورتی که آرشیو، دوره‌ای را پوشش دهد که طی آن بحران‌های بانکی توسعه یافته‌اند. به عبارت دیگر، به جای تحلیل یک وبگاه منفرد در سطحی خرد یا تحلیل همه وب در سطح کلان، با مجموعه‌های متمرکز وب در سطح میانی چه کار می‌توانیم بکنیم؟ قسمت عمده‌ای از پژوهش علوم اجتماعی در فضای برخط، در تعاملات سطوح میانی مشاهده می‌شود. آیا می‌توانیم برای فهم این که چگونه تغییر وب، واقعیت اجتماعی را منعکس و تقویت می‌کند یا تغییر می‌دهد، به صورت یکسان عمل کنیم؟

پاره ای از اقدامات برای پشتیبانی آرشیوهای وب موجود غیر ممکن خواهد بود - ممکن است برای انجام آن داده یا محتوا گردآوری نشده و تقریباً از بین رفته باشد. با این حال، رو به سوی تا آینده چه تغییراتی را می توانیم امروز و در سال های آینده برای آرشیوهای وب، ایجاد کنیم محققان در سالهای 2015، 2020 یا 2050 قادر به طراحی منابعی باشند که در حال حاضر برای پاسخ به این سؤال ها گردآوری می کنیم. محققان آینده، چه چیزی را از ما طلب می کنند که حالا در سال 2011، و در آینده، انجام نمی دهیم؟ مؤسسه های خاص چه می توانند انجام دهند؟ اگر IIPC به صورت جمعی در بهره داری از قدرت آرشیوهای چندگانه، عمل نماید چه اتفاق بهتر و مؤثرتری ممکن است رخ دهد؟

وب مجتمع: زندگی آرشیو وب

سؤال: چرا آرشیوهای وب به آرشیو شدن نیاز دارند؟ چرا آن ها نمی توانند با وب پویا یکپارچه گردند و به طور شفاف برای عموم و محققان قابل دسترسی باشند؟ امکان تصور وبی که با وب پویای فعلی بر روی صفحه ای قابل دسترسی به صورت منبعی پیش فرض از داده ها و اطلاعات لایه بندی شده باشد، وجود دارد. با این حال، این سطح می تواند بر روی لایه های زیرین وب قدیمی ایجاد شود که به آسانی برای علاقه مندان با حرکت نزولی یک لایه یا تعداد بیشتری از لایه ها به سمت پایین، قابل دسترسی است. اگر ده ها هزار ابزار قابل دسترسی برای پژوهش وب پویا بتواند برای این لایه های زیرین با استفاده از مکانیسم های ساده به کار گرفته شود احتمال کشف کاربردهایی برای اطلاعات و داده های موجود در وب گذشته توسط محققان افزایش می یابد.

احتمالاً این بزرگترین و بلند پروازانه ترین چالش در این گزارش است، زیرا نیاز به تغییر زیادی در زیر ساخت وب دارد. در عین حال، به معنی احتمال ضعیف وقوع آن است، اما نتایج مورد انتظار آن بسیار خواهد بود فراتر از ارزش پژوهشی که وب قدیمی به عنوان لایه های زیرین وب فعلی، دارد، به صورت بالقوه نیز باعث تغییر ساختاری در نحوه نگرش کاربران به وب خواهد بود. وب فعلی از نظر بسیاری برای اشاعه پیوندهای از دست رفته، از دست دادن اطلاعات، گم شدن صفحه ها، تغییر نشانی ها و تغییر اطلاعاتی که روی ویرایش های قبلی بدون هیچ امکانی برای مشاهده یا برگرداندن به ویرایش های قبلی، بازنویسی می شوند، غیر قابل اطمینان است. اگر ساختار اینترنت در برگشت زمان به عقب، به یکی از لایه های چندگانه تبدیل شود، به نحوی که حفره های موجود در لایه بالایی باعث ایجاد حفره در وب نشود، و به جای آن لایه پایین تر آشکار شود امکان اینکه وب به صورت منبعی با ثبات قابل اعتماد و مقاوم به از دست رفتن اطلاعات باشد، وجود خواهد داشت.

مشکل پیوندهای از دست رفته که تحلیل پیوند (1) نامیده می شود، مشکل دیرپایی برای کاربران وب پویاست. مشکل، زمانی پیچیده می شود که توجه بر وب آرشیوی متمرکز شود که روی هم رفته شناسه - گره های ثابتی برای ویرایش های آرشیوی صفحه ها وب ندارد. تلاش های متعددی صورت گرفته است،

به عنوان مثال (1) / (Web Cite (<http://www.webcitation.org>)) به نویسندگان امکان آرشیو نسخه ای از صفحه وب و ایجاد پیوند محافظت شده یا تجزیه کننده DOI را می دهد. (2) (Dead URL (<http://deadurl.com>)) رهیافت متفاوتی با اتکا بر آرشیو اینترنت و مخزن گوگل، در میان منابع دیگر برگزیده است تا سعی کند نسخه های پیوندهای از دست رفته را بیابد. اگر چه، تلاش های این چینی بوسیله اکثریت وسیعی از محققان که عمدتاً از وجود آن آگاه نیستند، بدون استفاده باقی می ماند. فراتر از عامل مشکل ساز، نتیجه ناخواسته و غیر منتظره دیگری نیز وجود دارد. عادت محققان به اینکه تا حد امکان نشانی های متعددی را در آثار علمی شان درج کنند. این امر اثرات متعددی دارد:

نخست: زمانی که سعی در ارزیابی تأثیر منابع برخط با استفاده از فنونی چون وب سنجی دارند، فقدان پیوندها باعث کاهش تأثیر منابع شان می شود.

دوم: خوانندگان را وادار به تلاشی جدی تر برای پیگیری منابع اطلاعاتی می کند، به طوری که سعی می کنند نه تنها منبع صحیح استناد، بلکه منبع استناد شده ای را که ممکن است ویرایش صفحه هایش می کنند به طور قابل توجهی تغییر کرده باشد نیز کشف کنند. اگر آرشیوهای وب تبدیل به منبعی قابل اتکا برای استناد پیوسته اطلاعات بشوند، این امر باعث افزایش پژوهش می شود و وجهه (نیمرخ) آرشیوهای وب را به صورت عمومی تر رشد خواهد داد.

در این وب مجتمع، دانش هویدا در وب، رشد و تکامل می یابد، اما به همان روش مشابه وب فعلی، به دور انداخته نمی شود. وب مجتمع، با استفاده از موتورهای جست و جویی نظیر گوگل، قابل خزش، زدودن پیوند شدن و تحلیل و جست و جو خواهد شد. پیوندها از بین نخواهند رفت، بلکه به سوی تبدیل شدن به و موادی هدایت می شوند. که مدت زیادی از وجودشان به صورت فعال بر روی وب فعلی نمی گذرد.

چالش بلند مدت: دو چالش در این سؤال نهفته است، که هر دوی آن ها مستلزم مشارکت کنندگان و فعالان زیادی خواهد بود. نخست: ما باید دوباره در مورد اینکه چگونه اینترنت را ببینیم و مهندسی کنیم فکر کنیم، در حال حرکت از موجودیتی تک لایه با پیوندهای جانبی بسیار به سوی موجودیتی چند لایه با پیوندهای جانبی فعلی، بلکه با پیوندهای فعلی به مواد قدیمی و پیوندهای قدیمی به مواد قدیمی یا فعلی [باشیم]. این چالش کم اهمیتی نیست و متقاعد کردن بسیاری از نقش آفرینان که در ساختن حال آینده اینترنت، سرمایه گذاری کرده اند، برای انتخاب، مشکل خواهد بود. با این حال، نتیجه آن منجر به زیر ساختی خواهد شد که وب گذشته را بسیار قابل دسترس تر برای ارجاع و جست و جو خواهد ساخت. چالش بزرگ بعدی این خواهد بود که آرشیویست های وب نیاز به تعریف مجدد نقش شان خواهند داشت، در حقیقت نه برای آرشیویست بودن به شکل کاملاً سنتی اش، بلکه به صورت متخصصانی که می توانند به پژوهشگران در درک روندها و منابع اینترنت در طول زمان کمک کند و در ابزارهایی که برای دسترسی و دستکاری لایه های این اینترنت چند لایه به آن ها نیاز دارند، برای راهنمایی پژوهشگران و عموم در پاسخ به پرسش های پیش بینی نشده در مورد رشد وب همانگونه که به صورت جدید توسعه یافته اند خبره هستند.

سؤال: محققان چگونه به تغییر رویدادها در جهان پاسخ می دهند و یا حوادث جاری را پیش می کنند؟ به طور فزاینده، رویدادهای محلی و جهانی خارج از وب اتفاق می افتد. این حوادث ممکن است ---ت رویدادهای مهم بین المللی و جالب توجه باشند. ممکن است مانند رویدادهای سیاسی اخیر در آفریقای شمالی و خاور میانه یا زلزله هایی در هائیتی ژاپن و سایر مکان ها رویدادهای محلی یا منطقه ای مهمی باشند. ممکن است کوچک اما گسترده، یا رویدادهای جاری باشند که در اولویت توجه گروهی کوچک یا حتی یک محقق باشند. این امر، به طور بالقوه انواع بینش ها را در مورد ماهیت و چیرستی اطلاعاتی که مردم به اشتراک می گذارند، موضوع ها و حوادثی که به طور برجسته توسعه می یابند، و چگونگی پاسخ افراد و سازمان ها و دولت ها به بحران ها و در طی زمان، چگونگی استحاله و افول حوادث در گفت و گوهای عامه مردم، به بار می آورد. این مورد ساده ترین چالش ولی واضح ترین آن هاست.

علاوه بر این، تعداد زیادی از محققان در این حوزه با هم کار می کنند، پدیده برداشت (1) موضوع چندین گفتگو در نشست IIPC در سال 2011 بود، توسط سخنرانانی که مثال های زیادی از پدیده برداشت در ارتباط با انقلابهای 2011 از سراسر دنیای عرب تا نشت نفت در آب های عمیق، تا المپیک سال 2012 لندن را بیان می داشتند.

یکی از موضوع های محوری کار با آرشیوهای وب، این است که ما نیاز به حرکت از فهم وب به صورت مجموعه ای از داده های انتخاب شده داریم و به جای آن باید آن را به صورت یک شبکه در حال تغییر و تحول ببینیم که به دوره های زمانی و رهیافت های طولانی نیاز دارد در این زمینه، تلاش هایی انجام گرفته است. برای نمونه، پروژه اروپایی تحلیل های طولی آرشیوهای وب (2) در حال ایجاد یک رصد خانه مجازی وب برای انجام تحلیل های طولی است.

چالش فوری: ایجاد سازوکارهایی برای محققان به منظور پیشنهاد سریع گرانولیتته (3) افزایشی و دامنه مناسبی برای استفاده در سایت های آرشیوی و موضوع های در حال تغییر، در حال حاضر، ممکن است یک پژوهشگر ماهر، ابزارهایی را برای تکرار اجرای خزش، راه اندازی کند، اما پژوهشگری با مهارت کمتر از نظر فنی بدون پشتیبانی قوی سازمانی، یادگیری نزولی خواهد داشت. زمانی که پدیده ای به سرعت متغیر، در حال توسعه است، پژوهشگر متخصص علاقه مند به آن پدیده که هیچ تجربه ای در وب آرشیوی ندارد، باید راهی برای گردآوری داده برای تحلیل آن، قبل از دست رفتن داشته باشد. سازمان های ماهر در راه اندازی ابزارهایی برای گردآوری داده های موجود بر روی این موقعیت های توسعه یافته، در سطح مناسبی از گرانولیتته، می توانند راه هایی برای محققان یا دیگران فراهم کنند که وبگاه ها، موضوع ها، کلید واژه ها، مانند آن را برای پاسخگویی سریع به رویدادهای وب برگزینند.

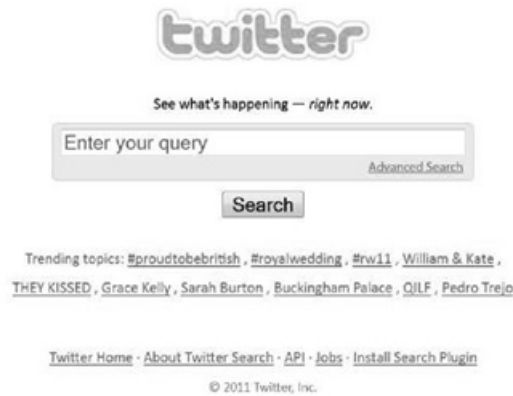
ص: 132

Harvesting -1

2- (http://www.lawa-roject.eu/the European Longitudinal Analytics of Web Archives project)

3- granularity (گرانولیتته، دانه دانه بودن) اصطلاحی در مورد اندازه واحد با مفهوم، استاندارد برای یک حالت عملیاتی خاص. اصطلاح مزبور در عناوینی چون وضوح صفحه نمایش معیار جست و جو و مرتب سازی پایگاه داده ها سطوح پردازش داده ها و مقدار زمانی که ریز پردازنده به یک عملیات اختصاص می دهد به کار می رود) (منبع: اصطلاحنامه تخصصی مایکروسافت).

آینده آرشیو وب ۱۳۳



تصویر ۱۹. روند موضوعات تویتر ، ۲۹ آوریل ۲۰۱۱

چالش در حال توسعه: استفاده از ابزارهایی نظیر خوراک‌های R.S.S. برای به‌کارگیری نشانه‌ها، بیانگر این است که تغییرات صفحه‌های وب نیاز به آرشیو شدن دارند.

در ارتباط با توسعه رویدادها، امکان پایش مواردی نظیر خوراک‌های R.S.S. یا نرم‌افزارهای کاربردی اخیراً توسعه یافته، برای داشتن سیستم‌هایی که به فعالیت تکثیر به وسیله افزایش تکرار ذخیره صفحه وب، پاسخ داده یا بوسیله آگاه‌سازی متصدیان انسانی حوزه‌های در حال توسعه بالقوه مورد علاقه، وجود دارد. چالش طولانی مدت: ساخت الگوریتم‌هایی که از روندهای فعالیت برخط (نظیر روندهای گوگل، یا روند موضوع‌های تویتر) به‌منظور راه‌اندازی گرانولیت‌آرشیوسازی افزایشی برای صفحه‌های وب مرتبط با آن موضوع‌های استفاده می‌کنند.

این امر مستلزم مهارت و پختگی بیشتری برای آرشیو‌هایی است که به‌طور مناسبی برای استفاده مجدد و به اشتراک گذاری و استانداردسازی ایجاد شده و استمرار دارند. پژوهشگران علاقه‌مند به استفاده از چنین آرشیو‌های گردآوری شده الگوریتم گونه‌ای (دارای الگوریتمی)، می‌خواهند منطبق حذف یا افزودن را بدانند و قادر به فهم ماهیت و محتوای مجموعه باشند. سؤال محوری که پرسیده خواهد شد این است که چگونه مجموعه‌ها با اکوسیستمی از منابع پژوهشی قابل استفاده و قابل دسترسی، متناسب شوند؟

کاربردهای آرشیوها و وبگاه‌ها

سؤال: چگونه افراد آرشیو‌های وب را بسیار مهم‌تر از وبگاه‌ها به‌کار می‌گیرند؟ در حال حاضر، احتمال زیادی وجود دارد که وبگاه‌های بر روی وب را در مقاطع مشخص زمانی، با استفاده از آرشیو‌های وب و زیر ساخت آرشیوی وب، مشاهده نمود. برای پژوهشگران دانشگاهی و صنعت تحلیل‌های

تصویر 19 روند موضوعات تویتر ، 29 آوریل 2011

چالش در حال توسعه: استفاده از ابزارهایی نظیر خوراک‌های R.S.S. برای به‌کارگیری نشانه‌ها بیان گر این است که تغییرات صفحه‌های وب نیاز به آرشیو شدن دارند.

در ارتباط با توسعه رویدادها، امکان پایش مواردی نظیر خوراک های R.S.S. یا نرم افزارهای کاربردی اخیراً توسعه یافته برای داشتن سیستم هایی که به فعالیت تکثیر به وسیله افزایش تکرار ذخیره صفحه، وب پاسخ داده یا بوسیله آگاه سازی متصدیان انسانی حوزه های در حال توسعه بالقوه مورد علاقه، وجود دارد.

چالش طولانی مدت: ساخت الگوریتم هایی که از روندهای فعالیت برخط (نظیر روندهای گوگل یا روند موضوع های توییتر) به منظور راه اندازی گرانولیتة آرشیو سازی افزایشی برای صفحه های وب مرتبط با آن موضوع های استفاده می کنند.

این امر مستلزم مهارت و پختگی بیش تری برای آرشیو هایی است که به طور مناسبی برای استفاده مجدد و به اشتراک گذاری و استاندارد سازی ایجاد شده و استمرار دارند پژوهش گران علاقه مند به استفاده از چنین آرشیوهای گردآوری شده الگوریتم گونه ای دارای الگوریتمی ، می خواهند منطق حذف یا افزودن را بدانند و قادر به فهم ماهیت و محتوای مجموعه باشند. سؤال محوری که پرسیده خواهد شد این است که چگونه مجموعه ها با اکوسیستمی از منابع پژوهشی قابل استفاده و قابل دسترسی متناسب شوند؟

کاربرد های آرشیو ها و وب گاه ها

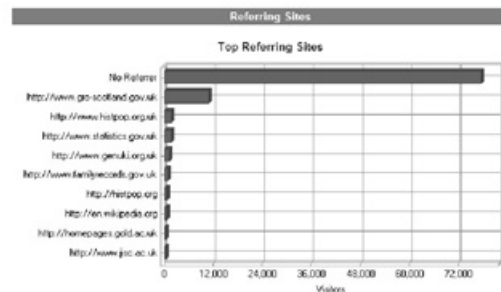
سؤال چگونه افراد آرشیوهای وب را بسیار مهم تر از وب گاه ها به کار می گیرند؟ در حال حاضر، احتمال زیادی وجود دارد که وب گاه های بر روی وب را در مقاطع مشخص، زمانی با استفاده از آرشیو های وب و زیر ساخت آرشیوی، وب مشاهده نمود برای پژوهش گران دانشگاهی و صنعت تحلیل های

سرویس دهنده ثبت وقایع و تحلیل های تحلیل گران تکنیکی پر تکرار برای ارزیابی کاربردها تأثیر و الگوهای ترافیک در وب پویاست با وجود، این فنون برای وب آرشیوی امکان پذیر نیست زیرا داده ها برای درک استفاده های وب قدیمی و آرشیو های وب به سادگی در دسترس نیستند.

عکس

۱۳۴ مدیریت منابع اطلاعاتی وب

سرویس دهنده ثبت وقایع و تحلیل های تحلیل گران، تکنیکی پر تکرار برای ارزیابی کاربردها، تأثیر و الگوهای ترافیک در وب پویاست. با وجود این، این فنون برای وب آرشیوی امکان پذیر نیست. زیرا داده ها برای درک استفاده های وب قدیمی و آرشیو های وب به سادگی در دسترس نیستند.



تصویر. ۲۰ نمونه ای از فایل داده ثبت وقایع برای سایت histpop.org منبع: مه پر و دیگران، ۲۰۰۹

چالش فوری: آرشیو سرویس دهنده ثبت وقایع سایت های وب آرشیوی، که محققان بتوانند نحوه استفاده از آرشیو های وب را مطالعه نمایند. این راه حلی ساده و سراسر است و در دسترس برای مؤسسه هایی است که آرشیو های وب را می سازند. سرویس دهنده ثبت وقایع آرشیو های وب می تواند برای پژوهشگران علاقمند به فهم چگونگی مرور آرشیو های وب توسط کاربران و چگونگی دسترسی آنها به منابع و اینکه چه بخش هایی از آرشیو بیشترین تکرار استفاده را دارند، ذخیره، نگهداری و قابل دسترسی شود. این اطلاعات به طور عمده مورد توجه جامعه آرشیوی وب خواهد بود. اما این گام نخست است.

چالش بلند مدت: تلاش بلند پروازانه تری است، اما علاقه بالقوه بسیار گسترده تری خواهد بود: راه اندازی زیر ساختی برای امکان آرشیوسازی سرویس دهنده های ثبت وقایع و تحلیل های مرتبط و پیوند با وبگاه ها، به طوری که پژوهشگران نه فقط آنچه که بر روی وب موجود بوده، بلکه نحوه استفاده از آن را نیز بتوانند ببینند. این هدفی بسیار بلند پروازانه تر است، زیرا سرویس دهنده های ثبت وقایع و اعتبار های تحلیل ها فقط به صورت داخلی برای سرور و مدیران اعتبار در وضعیتی حفاظت شده قابل مشاهده هستند. تجربیات عمومی مدیران سرور برای ذخیره سرویس دهنده ثبت وقایع در بلند مدت ضروری نیستند، به طوری که آنها به طور عادی برای ذخیره سازی فضا و دوری از انباشتگی و آشفستگی سرور، حذف و بازنویسی می شوند. اگر چه این داده ها، به طور بالقوه برای پژوهشگرانی ارزشمند است که نه فقط می خواهند از وضعیت یک سایت آگاه شوند، بلکه می خواهند از چگونگی و مقدار استفاده از آن، ترافیک منابع و دیگر حقایقی که از ثبت وقایع و داده های تحلیلی می تواند گردآوری شود، آگاهی یابند. راه حل های اجتماعی شامل ایجاد مکانیسم هایی برای مدیران سرور برای مشارکت در ثبت وقایع

تصویر. 20 نمونه ای از فایل داده ثبت وقایع برای سایت histpop.org منبع: مه پر و دیگران، 2009

چالش فوری: آرشیو سرویس دهنده ثبت وقایع سایت های وب آرشیوی، که محققان بتوانند نحوه استفاده از آرشیوهای وب را مطالعه نمایند این راه حلی ساده و سر راست و در دسترس برای مؤسسه هایی است که آرشیوهای وب را می سازند. سرویس دهنده ثبت وقایع آرشیوهای وب می تواند رای پژوهش گران علاقمند به فهم چگونگی مرور آرشیوهای وب توسط کاربران و چگونگی دسترسی آن ها به منابع و اینکه چه بخش هایی از آرشیو بیش ترین تکرار استفاده را دارند، ذخیره، نگهداری و قابل دسترسی شود. این اطلاعات به طور عمده مورد توجه جامعه آرشیوی وب خواهد بود. اما این گام نخست است.

چالش بلند مدت: تلاش بلند پروازانه تری، است اما علاقه بالقوه بسیار گسترده تری خواهد بود راه اندازی زیر ساختی برای امکان آرشیو سازی سرویس دهنده های ثبت وقایع و تحلیل های مرتبط و پیوند با وب گاه ها به طوری که پژوهش گران نه فقط آن چه که بر روی وب موجود بوده، بلکه نحوه استفاده از آن را نیز بتوانند بینند این هدفی بسیار بلند پروازانه تر است، زیرا سرویس دهنده های ثبت وقایع و اعتبارهای تحلیل ها فقط به صورت داخلی برای سرور و مدیران اعتبار در وضعیتی حفاظت شده قابل مشاهده هستند. تجربیات عمومی مدیران سرور برای ذخیره سرویس دهنده ثبت وقایع در بلند مدت ضروری نیستند به طوری که آن ها به طور عادی برای ذخیره سازی فضا و دوری از انباشتگی و آشفستگی سرور، حذف و بازنویسی می شوند. اگر چه این داده ها به طور بالقوه برای پژوهش گرانی ارزشمند است که نه فقط می خواهند از وضعیت یک سایت آگاه شوند، بلکه می خواهند از چگونگی و مقدار استفاده از آن، ترافیک منابع و دیگر حقایقی که از ثبت وقایع و داده های تحلیلی می تواند گردآوری شود، آگاهی یابند. راه حل های اجتماعی شامل ایجاد مکانیسم هایی برای مدیران سرور برای مشارکت در ثبت وقایع

است که می تواند با وبگاه های آرشیو ارتباط داشته باشد و برای ارائه کنندگان تحلیل هایی همچون گوگل به منظور ارائه گزینه ای به منظور مشارکت سایت های تحلیلی با آرشیوهای وب و تا حد امکان تعیین دوره محدودیت قبل از این که داده ها انتشار یابند.

چالش بلند پروازانه طراحی نظام هایی برای آرشیو نه تنها سرویس دهنده های ثبت وقایع، بلکه خود ترافیک وب در یک سبک ناشناخته حفاظت شده این راه حل حتی بسیار بلند پروازانه تر است؛ زیرا سازو کارهای تحلیل ترافیک، وب روی هم رفته عموماً در دسترس نیستند. نگرانی های نهفته ای از مواردی نظیر بازرسی بسته های عمیق که به تحلیلگران درباره ترافیک روان، وب می گوید وجود دارد.

پرسش محوری برای سؤال کردن این است که چه زمان مزایای فهم نحوه رفتار افراد در وب بر خطرات غلبه مینماید، بنابراین توجه به این امر که روش هایی وجود دارند که این داده ها می توانند برای تحلیل های بعدی آرشیو و ذخیره شوند در زمانی که خطرات برای افراد و سازمان ها به حد کافی به واسطه گذر زمان و گمنام کردن داده ها کاهش می یابد ارزشمند است.

متخصصی وب

سؤال: آیا امکان شناسایی مجموعه هایی که اندازه و شکل وب تاریخی را در رابطه با حوزه های تخصصی مورد توجه می سنجند، وجود دارد؟ فرض می کنیم که بسیاری از گروه ها و تنالگان ها وب گاه هایی را در طول دوره های زمانی ایجاد خواهند کرد چگونه امکان دارد که شناسایی تداوم حضور آن ها در وب و جمع آوری بدنه مناسبی از سایت های پژوهشی، برای آن ها کوچک و بی اهمیت باشد؟

مثال های متعددی را می توان بر شمرد گروه های علاقه مند به سرگرمی موضوع های تخصصی دانشگاهی میراث صنایع دستی نظیر دیداری و شنیداری سایت های گروه های سیاسی، و مانند آن. مطالعه این موارد می تواند موجب ساخت ابر مجموعه ها شود مجموعه مجموعه ها، آرشیو - در - جعبه چه راهنمایی در زمان ساختن مجموعه های آرشیوی از منابع متنوع وب برای افزایش ارزش این مجموعه ها مورد نیاز است؟ یک ابر مجموعه تا چه اندازه می تواند با دیگری مقایسه و پیوند داده شود و حتی به بخش هایی از ابر مجموعه های بزرگ تر تبدیل شوند؟

متخصص وب، ضرورت مقیاس کوچک تری از داده های انتخاب شده بر اساس رویداد ها یا موضوع هایی هم چنان، انتقادی و موازی با تجربیات و انتظارات موجود پژوهشی را تشخیص می دهد. مجموعه قوانینی از این دست به صورتی قابل مشاهده و جست و جو بوسیله افراد یا تیمی باقی میماند که این رشته های داده را از یک چشم اندازه ویژه رشته یا موضوع پژوهشی مهمی در دنیای آکادمیک، ایجاد و تحلیل می کنند به ویژه در زمینه هایی که باوری دیرینه وجود دارد که پژوهش گران، خالقان مدیران و تحلیل گران مجموعه قوانین خود شان باشند با این حال حتی این مجموعه های تخصصی به طور بالقوه در زمانی که به طور استاندارد برای پاسخ به سؤال های محققان دوباره ترکیب و دسته بندی می شوند، می توانند ارزش آفرین باشند

این، امر مسئله مقیاس پذیری را در مسیر پرسش های مرتبط با آرشیو های وب متخصص پرورش

می دهد. آیا اندازه بحرانی در طول دوره های پوشش و دامنه مجموعه وب آرشیوی برای اینکه در طول زمان مورد استفاده دیگر پژوهش گران باشد وجود دارد؟ چگونه راهبرد ها و سیاست های خزندگی وب مؤسسه های مختلف (همچون برداشت های عمده و حوزه رویدادی یا انتخابی) انتظارات پژوهش گران را برآورده کرده، و چه راهبرد های متنوعی به مناسب ترین وضع برای پشتیبانی، استفاده می شوند؟

بعضی از مثال های ویژه حوزه مجموعه های تخصصی در ابزارهای علوم اجتماعی برای ترکیب بانک جهانی داده با تحلیل های موجود در ولف رام آلفا (1) در ارتباط با سلامتی یا دیگر اطلاعات جمعیتی می توانند توسعه یابند. این مستلزم استفاده از ابزاری (همچون ولف رام آلفا) دارد که به طور مداوم با داده های جدید روزآمد می شود به همان صورت الگوریتم ها و ابزار های مصور سازی آن، بستر تحلیلی داده های پویا را ارائه می کنند در علوم، طبیعی پژوهش گران به گرفتن اطلاعات آب و هوا (دما) علاقه مند خواهند بود و این اطلاعات جمع آوری شده را با طبیعت شناسان غیر حرفه ای (مانند مشاهده کنندگان پرندگان) در گروه های مختلف (با استفاده از ابزارهایی چون گوگل پلاس یا فیس بوک) در سراسر دنیا، به اشتراک گذاشته و به گروه هایی می پیوندند که چگونگی تغییر الگوهای مهاجرت پرندگان با آب و هوا یا آن چه مهاجرت پرندگان درباره تغییرات آب و هوا می تواند بگوید را تحلیل می نمایند. این امر نیاز به آشنایی با علم، شهری جامعه ای برخط و محیط گرایی دارد در علوم انسانی به عنوان عقاید تاریخی ممکن است مصاحبه های آلن مک فارلن (2) با متفکران معاصر اصلی که در سخنرانی های (3) iTunesu آکسفورد نمایش داده می شود مقایسه شود تا الگوهای چگونگی تفکر اجتماعی متفکران مهم در مقابل آن هایی که در iTunesu مشهور هستند شفاف شود. (به عبارت دیگر مقایسه افکار اصلی در مقابل با افکار مردمی).

چالش فوری: ارائه راهنمایی برای انواع متفاوتی از ابر مجموعه ها به طوری که عناصر استاندارد در آرشیو وب وارد شوند و در نتیجه دامنه و انسجام آن ها تضمین و تعیین شود.

چالش توسعه یافته: ارائه ابزارها و زیر ساخت های سازمانی برای این که پژوهش گران بتوانند بر تردیدها غلبه یابند تا ابر مجموعه هایی را بسازند که حداقل به صورت بالقوه توسط یک پژوهش گر منفرد مورد استفاده قرار گیرد. استانداردهایی تعریف شود تا ابر مجموعه های آرشیوی توسط یکدیگر قابل استفاده باشند.

چالش بلند مدت: تشویق به پدید آمدن سازمان هایی که ابر مجموعه ها را به طور مفید و گسترده تشویق و حمایت می کنند.

وب دیداری

سؤال: اگر بخواهم از تصاویر برای فهم چگونگی تغییرات دنیا استفاده کنم آیا می توانم برای فهم این فرآیندها به صورت دیداری تصاویر را از آرشیوهای وب استخراج کنم؟ برای مثال: آیا این امکان وجود دارد که با استخراج تصاویر تغییر یافته در طول زمان از مکان های یکسان بر روی صفحه هایی نظیر

ص: 136

(1) [Flicker](#) تحلیل های تصویری انجام گیرد؟ تصویر برداری مجدد، تجربه ای است که طی آن مکانی که قبلاً از آن عکس برداری شده بازدید کرده و برای مستند کردن تداوم و تغییرات در طول زمان، عکس جدیدی از آن گرفته می شود. یکی از نخستین پروژه ها برای انجام این کار بازدید مجدد از 1200 سایت در آمریکای غربی و تصویر برداری از صد سال قبل به وسیله پژوهشگران دولتی بوده است. (کلت، (2) منچستر (3) و وریورگ 1984. (4)).

حالا، تصور کنید 100 سال بعد از حال را که قادر باشیم نه تنها عکس های یکسانی را از مکانی مشخص مقایسه کنیم بلکه بتوانیم از آرشیو های وب سری کامل عکس های مستند شده ای را از دنیای در حال تغییر و ثابت در طول زمان استخراج کنیم.

عکس

Flicker^۱، تحلیل های تصویری انجام گیرد؟ تصویر برداری مجدد، تجربه ای است که طی آن مکانی که قبلاً از آن عکس برداری شده، بازدید کرده و برای مستند کردن تداوم و تغییرات در طول زمان، عکس جدیدی از آن گرفته می شود. یکی از نخستین پروژه ها برای انجام این کار، بازدید مجدد از ۱۲۰۰ سایت در آمریکای غربی و تصویر برداری از صد سال قبل به وسیله پژوهشگران دولتی بوده است. (کلت، منچستر و وربورگ، ۱۹۸۴)

حالا، تصور کنید ۱۰۰ سال بعد از حال را که قادر باشیم نه تنها عکس های یکسانی را از مکانی مشخص مقایسه کنیم، بلکه بتوانیم از آرشیوهای وب سری کامل عکس های مستند شده ای را از دنیای در حال تغییر و ثابت در طول زمان استخراج کنیم.



تصویر ۲۱. تصویر قصر باکینگهام که از ۲۲۰ تصویر مختلف گردآوری شده است.

منبع: <http://photosynth.net/view.aspx?cid=34e49d3e-2d1e-4118-bbad-d2f5d74ce340>

چالش فوری: تضمین اینکه، تصاویری که بیشتر در صفحه های وب آرشیو شده، از دست می روند، در اولویت حفاظت قرار گیرند.

چالش توسعه یافته: ساخت فناوری هایی همانند photosynth^۲، که قادر به اتصال تعداد زیادی از تصاویر برای نمایش پانورامای یک مکان یا شی برای کار با اطلاعات زمانی به منظور جمع آوری مناظر مشابه در طول زمان باشند.

چالش بلند مدت: ایجاد یک آرشیو از تصاویر جهان، شامل اطلاعات بسیار از زمان و مکان که تا حد امکان برگرفته از EXIF^۳ و داده های صفحه های وب باشند، به طوری که تصاویر بتوانند برای پژوهش استفاده شوند. ابزارهایی برای جای دادن، استخراج، ترکیب و تصاویر مورد نیاز است.

1. Flickr
2. Klett
3. Manchester
4. Verburg
5. <http://photosynth.net/Background.aspx>.
6. EXIF

تصویر 21. تصویر قصر باکینگهام که از 220 تصویر مختلف گردآوری شده است.

منبع: <http://photosynth.net/view.aspx?cid=34e49d3e-2d1e-4118-bbad-d2f5d74ce>

چالش فوری: تضمین، این که تصاویری که بیش تر در صفحه های وب آرشیو شده، از دست می روند، در اولویت حفاظت قرار گیرند.

چالش توسعه یافته ساخت فناوری هایی همانند (5) <http://photosynth.net/Background.aspx>. که قادر به اتصال تعداد زیادی از

تصاویر برای نمایش پانورامای یک مکان یا شی برای کار با اطلاعات زمانی به منظور جمع آوری مناظر مشابه در طول زمان باشند.

چالش بلند مدت: ایجاد یک آرشیو از تصاویر جهان شامل اطلاعات بسیار از زمان و مکان که تا حد امکان برگرفته از [EXIF 6](#) و داده های صفحه های وب باشند به طوری که تصاویر بتوانند برای پژوهش استفاده شوند. ابزارهایی برای جای دادن استخراج ترکیب و تصاویر مورد نیاز است.

ص: 137

Flicker -1

Klett -2

Manchester -3

Verburg -4

photosynth -5

EXIF -6

وب همان گونه که بود

سؤال: چگونه می توانم وب را همان گونه که بود ببینم؟ اگر من بخواهم وب را همان گونه که، در یکم ژانویه 2011 بود مرور کنم و قادر باشم تا روی صفحه ها، تصاویر پیوندها و دیگر محتوای آن همان گونه که در آن روز ظاهر شده بود کلیک کنم چگونه این کار را می توانم انجام دهم؟

عکس

وب همان‌گونه که بود

سؤال: چگونه من می‌توانم وب را همان‌گونه که بود ببینم؟ اگر من بخواهم وب را همان‌گونه که، در یکم ژانویه ۲۰۱۱، بود مرور کنم و قادر باشم تا روی صفحه‌ها، تصاویر، پیوندها و دیگر محتوای آن همان‌گونه که در آن روز ظاهر شده بود، کلیک کنم، چگونه این کار را می‌توانم انجام دهم؟



تصویر ۲۲- ویرایش بتای پاسخ WayBack Machine

(منبع <http://replay.web.archive.org/20041010185532/http://netpreserve.org/about/index.php>)

ویرایش بتای فعلی نسخه پاسخ [پایگاه] Wayback Machin چنین کارکردی را نوید می‌دهد (گشت و گذار در وب همان‌گونه که بود، نسخه بتا، که در حال حاضر در سایت تبلیغ می‌شود)، اما IIPC یا آرشیوهای انفرادی چه تلاش‌های دیگری می‌توانند برای امکان تکرارپذیری وب، انجام دهند؟

چالش: گسترش و افزایش تلاش‌هایی برای ایجاد وب فعلی به صورت تکرارپذیر در آینده که نیاز به گردآوری، ذخیره‌سازی و اشاعه مجدد وب فعلی همان‌گونه که در گذشته بود، خواهد داشت. سؤال محوری برای پرسش در اینجا این است که [وب] چگونه می‌تواند فراتر از این که منبعی ارجاعی صرفاً برای رفع کنجکاوی یا مراجعه گاه به گاه باشد، مورد استفاده قرار گیرد. چه نوع نیاز بکر یا سؤال پژوهشی غیر متصور بر جست‌وجوهای دستی و گشت و گذار در وب گذشته (قدیمی) اتکا دارد؟ آیا مورخان آینده، همان‌گونه که در دنیای امروز خواندن اخبار و انتشارات و دوران گذشته فانی وب مورد علاقه بوده‌است، به آن علاقه‌مند هستند؟ آیا آنها می‌خواهند از نرم‌افزارهای کاربردی استفاده کنند یا صرفاً از منابع مرجع؟ به عبارت دیگر، بزرگ‌ترین سؤال، ساختن موارد استفاده برای وب تکرارپذیر و سپس ساختن رابط‌هایی که این موارد استفاده را پشتیبانی می‌کند، می‌باشد. انجام این امر مستلزم مشورت متخصصان آرشیو سازی وب با متخصصان حوزه‌هایی شامل مورخان و دیگر افرادی است که به بازسازی مجدد گذشته علاقه‌مند هستند.

تصویر 22- ویرایش بتای پاسخ WayBack Machine

(منبع <http://replay.web.archive.org/20041010185532/http://netpreserve.org/about/index.php>)

ویرایش بتای فعلی نسخه پاسخ [پایگاه] Wayback Machin چنین کارکردی را نوید می‌دهد (گشت و گذار در وب همان‌گونه که بود، نسخه بتا که در حال حاضر در سایت تبلیغ می‌شود) اما IIPC یا آرشیوهای انفرادی چه تلاش‌هایی می‌توانند برای امکان تکرارپذیری وب، انجام دهند؟

چالش: گسترش و افزایش تلاش هایی برای ایجاد وب فعلی به صورت تکرار پذیر در آینده که نیاز به گردآوری، ذخیره سازی و اشاعه مجدد وب فعلی همان گونه که در گذشته بود، خواهد داشت. سؤال محوری برای پرسش در این جا این است که [وب] چگونه می تواند فراتر از این که منبعی ارجاعی صرفاً برای رفع کنجکاوی یا مراجعه گاه به گاه باشد مورد استفاده قرار گیرد چه نوع نیاز بکر یا سؤال پژوهشی غیر متصور بر جست و جوی دستی و گشت و گذار در وب گذشته (قدیمی) اتکا دارد؟ آیا مورخان آینده، همان گونه که در دنیای امروز خواندن اخبار و انتشارات و دوران گذشته فانی وب مورد علاقه بوده است، به آن علاقه مند هستند؟ آیا آن ها می خواهند از نرم افزارهای کاربردی استفاده یا صرفاً از منابع مرجع؟ به عبارت دیگر بزرگ ترین سؤال ساختن موارد استفاده برای وب تکرار پذیر و سپس ساختن رابطه ای که این موارد استفاده را پشتیبانی می کند می باشد. انجام این امر مستلزم مشورت متخصصان آرشیو سازی وب با متخصصان حوزه هایی شامل مورخان و دیگر افرادی است که به بازسازی مجدد گذشته علاقه مند هستند

سؤال: وب چگونه مقایسه شده و در طول زمان تغییر می کند؟

افزایش تلاش هایی برای فهم وب به صورت یک سیستم مستلزم قابلیت های افزایش به مقیاس بزرگ است که قادر خواهد بود الگوها و روندها را در طی زمان دست بی اندازد (تغییر دهد). برای انجام این امر، ما نیاز داریم که پرسیم چه رهیافت هایی برای توسعه روش های تحلیلی معتبر در دسترس هستند؟ چگونه ما می توانیم فرضیات ساخته شده درباره داده های وب را به صورت یک سری داده اعتبار سنجی کنیم؟ چه ابزارهای آماری برای مجموعه های وب آرشیوی می تواند به کار گرفته شود و چه ابزارهایی نیاز به توسعه دارند؟

حتی آمارهای ساده نیز برای استخراج درباره وب بی اهمیت نیستند برای نمونه، چه تعداد از وب گاه ها به طور سالانه (سراسر جهان در یک کشور مشخص روی موضوع مشخص) برای X تعداد سال گذشته بوده اند؟

در مجموعه آرشیوی (آرشیو در جعبه) تاریخ ایجاد صفحه ها چیست؟ صفحه های آن به چه زبان هایی هستند؟ آیا در زمان ایجاد صفحه ها روندهایی وجود داشته است؟ آیا خوشه ای است؟ آیا یک فرایند ساختمانی ثابت بوده است؟ آیا موضوع های مشخصی بیش تر از دیگران پیوند داده شده اند؟ آیا برخی از انواع مجموعه ها احتمال کم تر یا بیش تری برای پیوند به منابع خارجی دارند؟ آیا وب گاه ها می توانند به رده هایی که می توانیم با استفاده از تحلیل های خوشه ای کشف کنیم تقسیم بندی شوند؟

آیا ما می توانیم سایت ها را به وسیله آمار هایی همچون اندازه میانگین وب گاه ها در رده های مختلف، میانگین تعداد پیوند ها میزان داده های غیر متنی (عکس ها، تصاویر و مانند آن) سن محتوا در سن محتوا در فاصله بین روز آمدسازی ها دفعات روز آمدسازی، نوع رابط (ثابت در مقابل دینامیک به طور مثال)، مقایسه کنیم؟

چگونه این آمارها به ما در فهم ساختار مجموعه ها و وب کمک می کنند؟

چالش: ایجاد ابزارها و روش هایی برای استفاده از وب به عنوان یک رشته داده عظیم به جای مجموعه ای از اسناد در حال حاضر در صورت وجود داده ها اگر شخص بخواهد بداند چه چیزی باید به طور بسیار اساسی درباره اندازه و ساختار وب فعلی یا وب به صورت گذشته اش مورد سؤال باشد داده ها برای پژوهش گران قابل دسترسی نیستند. بنابراین ابزارهایی باید ایجاد شوند که آمارگیری روی وب یا روی صفحه های یک مجموعه آرشیوی را انجام دهند.

ایده ها چگونه تکثیر می شوند

سؤال: ایده ها چگونه روی اینترنت کشش یافته و تکثیر می شوند؟ یکی از جنبه های قابل توجه اینترنت توانایی شگفت انگیز آن برای پشتیبانی انتقال از الگوی رفتاری ایده هایی که رشد می یابند و گسترش فرهنگی است اگر علاقه مند هستیم به چگونگی ویروسی شدن ویدئو ها یا چگونگی گسترش شوخی ها یا چگونگی وارد شدن یک بیت از اطلاعات یا اطلاعات غلط به آگاهی عمومی، چه ابزارهایی به ما کمک خواهد نمود؟ چگونه ابزارهایی را با توانایی امکان این که یک آرشیو ساخته شود نه بر اساس

جغرافیایی فیزیکی یا مجازی بلکه بر اساس حرکت از یک ایده می‌سازیم؟ هر کس می‌تواند تصور کند که قادر به تعیین یک ایده باشد و برای دنبال نمودن آن ایده به صورتی که در طول زمان توسعه می‌یابد، حرکت کند

هم چنین چه زمینه گسترده‌ای در اطراف محتوایی که ما در آرشیو می‌بینیم، وجود دارد؟ برای مثال: مردم روی وب در زمان پیدایش آن چه چیزی را جست و جو می‌کردند؟ گوگل زیتگیست (1) و گوگل ترندز (2) در مورد چیزهایی که مردم درباره آن جست و جو می‌کردند به ما می‌گویند. چه چیزهایی دیگری را می‌توانیم برای فهم زمینه جمع‌آوری کنیم؟ برای مثال- در تویتر اشاره شده است - به فهم زمینه محتوا بوسیله مشاهده چیزهایی که با یکدیگر آمده‌اند، کمک می‌کند.

برای نمونه واتسون آی بی ام (3)، نظام اختصاصی است که از مواد آرشیوی زیادی برای ساخت موتور دیپ کیو ای (4) خودش استفاده می‌کند که به آن برای برنده شدن در بحران 2011 کمک می‌کند. چگونه این نوع از ابزار به طور گسترده تری برای جست و جوی پیشرفته در دسترس خواهند شد؟

چالش: بعد زمان وب همان گونه که ایجاد شده نیاز دارد که حفاظت شده قابل استخراج و قابل تحلیل، باشد گرانیته بهتری نیاز است تا بیند ایده‌ها از کجا آغاز شده، چگونه گسترش می‌یابند و این که چه فعالیت‌هایی سرعت تکثیر آن‌ها را کاهش یا افزایش می‌دهد؟ ایده‌ها به صورت موجودی زنده در دنیا پدید می‌آیند و تنها پس از این که نگهداری، شوند برای برگشت به اصلیت شان مورد توجه قرار خواهند گرفت با این حال بدون گرانیته و عمق مناسب آرشیوی ایده اصلی ممکن است در زمانی که شخص به جست و جوی آن فکر می‌کند، از دست رفته باشد.

وب غیر قانونی

سؤال: چگونه وب برای پشتیبانی و توانایی فعالیت‌های غیر قانونی استفاده می‌شود و چگونه در طول زمان تغییر می‌کند؟ نوعی محتوا که به صورت برخط تکثیر می‌شود و کم‌تر مورد توجه دانشمندان است، مواد غیر قانونی وب است این طیف گسترده‌ای از محتوای جنسی تا اطلاعاتی درباره استفاده از مواد، مخدر قمار منابع گروه‌های تندرو محتوای مرتبط با تروریسم و دیگر منابعی است که یا غیرقانونی یا دارای مشکلات اجتماعی هستند سؤال این جاست که چه کسی باید محتواهای غیر قانونی یا قانونی ولی کم‌تر از نظر اجتماعی پذیرفته شده وب را آرشیو کند؟ چگونه بدون شکستن قانون می‌توانیم این کار را انجام دهیم؟ و چگونه می‌توان برای محققان بدون در خطر افتادن هم پژوهش‌گر و هم مؤسسه‌هایی که دسترسی را ایجاد ساخته‌اند قابل دسترسی شود؟

دانستن این که چه فعالیت‌های غیر قانونی‌ای از نظر حجم و عمومیت یافتن در حال رشد هستند کدام یک در طول زمان رنگ می‌بازند و چه فعالیت‌های غیرقانونی غیر منتظره‌ای پدید می‌آیند نه فقط

ص: 140

برای پژوهش گران، بلکه برای سیاست گذاران عمومی، متخصصان حوزه سلامت که احتیاج به پیگیری نتایج رفتارهای خطر آفرین دارند خبرگان سلامت عمومی، متخصصان و مؤسسه های حمایت اجتماعی و کسانی که مسئول حفظ رفاه قشر آسیب پذیر هستند مفید می باشد

چالش: بزرگ ترین چالش در این جا این است که حتی اگر منابع غیر قانونی بر روی اینترنت رواج داشته باشند سازمان های اندکی هستند که خارج از اعمال فشار قانونی مایل به پذیرش خطر مشارکت در جمع آوری داده های مربوط به این منابع هستند. هستند تابوهای فرهنگی و خطرات قانونی در ارتباط با دسترسی و ذخیره سازی منابع غیر قانونی حتی برای اهداف مثبت همانند پژوهش فهم جنبه های جامعه مدرن بیشتر پژوهش گران وب و آرشیویست های وب را از چنین منابعی بر حذر می دارد.

به نظر می رسد مهم ترین سازوکاری که می تواند در این جا وجود داشته باشد نظامی با حفاظت قانونی خواهد بود که به طور مناسب و تا حد امکان افراد و سازمان ها را برای آرشیو و جست و جوی داده های غیرقانونی قابل دسترسی اینترنت بدون در معرض خطر قرار گرفتن سازمان ها یا پژوهش گران استفاده کننده از این مجموعه ها تأیید کند

رد پای رقومی

سؤال: چگونه می توانیم (و باید) رد پای رقومی یک شخص را آرشیو کنیم؟ فعالیت ها و اقدامات یک شخص به طور پیوسته مورد توجه بالقوه است به ویژه اگر شخص (مشهور باشد یا بشود) قبلاً در این مورد بحث شده است که (گارفینکل (1) و کاکس 2009 (2) فهرست آثار زندگی یک شخص برای آرشیویست ها یک وظیفه خواهد بود آرشیو وب یک شخص می تواند شامل صفحه های وب پروفایل های شبکه اجتماعی و پست هایش، ارتباطاتش، انتشارات، و منابع دیگر در مورد زندگی رقومی او باشد.

چالش: چالش محوری، فهمیدن چگونگی ساختن ابزارهایی است که به افراد اجازه می دهد تا به صورت دستی مشخص کنند که چگونه به صورت خودکار آثار رقومی خود را جمع آوری کنند. آیا ابزارها به طور خودکار ردپای رقومی را تا حد امکان بر مبنای انتخاب، جمع آوری می کنند؟ چگونه می توانیم نظام هایی را ایجاد کنیم که نه فقط یادآوری بلکه امکان فراموشی به وسیله حذف های بعدی توسط افراد (و فراموشی) بخش ها یا همه رد پاها را داشته باشد به صورتی که برخی دانشمندان از آن به عنوان یک حق اساسی ذکر کرده اند مه یر (3) شونبرگر 2009 (4)؟ پیشرفت ها به ویژه در این حوزه مشکل است به دلیل این که مسائل بسیار زیاد حقوقی و خصوصی آن نیاز به بررسی خواهد داشت.

ص: 141

Garfinkel -1

Cox -2

Meyer -3

Schonberger -4

سؤال: چگونه داده‌ها از آرشیوهای وب دوباره استخراج می‌شوند؟ در سال‌های اخیر، رشد چشم‌گیری در وب مبتنی بر داده (در مقابل وب اسنادی) صورت گرفته است مسائل متعددی مطرح شده است، مانند این که چگونه داده‌ها در کنار اسناد آرشیو بشوند؟ چه نوع ابزار پاک‌کننده داده‌ها برای کار با داده‌ها مورد نیاز خواهد بود؟

برای مثال به فرآیندهای بین حوزه‌های علمی یا صنعتی برای فراهم‌آوری دانش مانند طراحی، هوافضا، کشف مواد مخدر و غیره، فکر کنید زمانی که یک هواپیما دچار سانحه می‌شود و بازرسان می‌خواهند محاسبات اولیه مهندسی را دوباره ارزیابی کنند چگونه می‌توانیم توانایی درک داده‌ها از یک زنجیره تأمین طراحی مهندسی هفتاد ساله را حفظ کنیم؟ طراحی رقومی، بود دانش توسط 100 نفر همکار تدارک می‌شد که بعضی از آن‌ها در فواصل زمانی از حوزه کار خارج شده‌اند و همه آن توسط مجموعه کاملاً متفاوتی از افراد به کار گرفته می‌شد در این مورد چرخه زندگی دانش بسیار طولانی‌تر از چرخه زندگی کسب و کار است و آرشیوها نقش اساسی در حفظ این اطلاعات ایفا می‌نمایند.

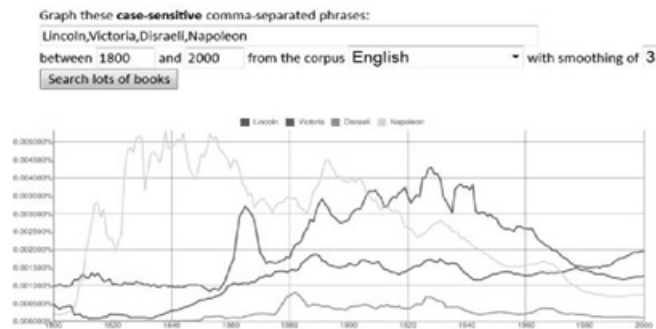
سؤال‌های مرتبط این است که چگونه می‌توانیم اطلاعات اختصاصی را آرشیو و تحلیل کنیم؟ بعضی از داده‌ها اختصاصی هستند و لازم است برای تعهد به موافقت نامه‌های حقوقی نگهداری شوند. با وجود، این اگر نهادهای قابل، اعتماد داده‌های خام را ذخیره کنند و فقط به ابزارهای تحلیلی اجازه دسترسی به آن‌ها را بدهند تحلیل‌های بسیاری بدون دسترسی به داده‌های خام می‌توانند انجام گیرند. سپس نتایج گردآوری و تلخیص شده می‌توانند بدون افشای داده‌های حفاظت شده در دسترس پژوهش‌گر قرار گیرند. ابزارهایی برای سرویس دهی به عنوان الگویی که امکان دسترسی به بخش‌هایی از داده‌های خام را می‌دهند به جای اینکه داده‌های خام را بارگذاری کنید وجود دارند همانند الیکسر، (1) برنامه‌ای برای دسترسی به اطلاعات نجوم مثال، دیگر مشاهده گر n-gram کتاب‌های گوگل (2) که به کاربران امکان تحلیل داده بدون دسترسی به داده خام را می‌دهد.

اگر یکبار به داده‌هایی از مخازن داده‌های قابل، تحلیل دسترسی داشته باشید، امکانات افزوده فزونی خواهند یافت. همانند پیوند داده‌ها از طریق، ابزارها ایجاد کتابخانه‌هایی متشکل از اجزایی که به پژوهش‌گران اجازه می‌دهند که داده‌ها را از راه‌های گوناگون تحلیل کنند و امکاناتی برای ترکیب و تطابق آن‌ها را به روش‌های جدید فراهم می‌کنند. اگر در ایجاد این سری داده‌ها و ابزارها، دقت شود، الگوها و همبستگی‌هایی که قبلاً کشف نشده‌اند، افزایش می‌یابد.

ص: 142

آینده آرشیو وب ۱۴۳

Google labs Books Ngram Viewer



تصویر ۲۳- مشاهدهگر n-gram کتابهای گوگل

چالش: برای بخش‌هایی از اینترنت که حاوی داده‌ها به جای اسناد هستند، تغییر نگرش از اینکه آرشیوهای وب برای نگهداری اسناد هستند به سمت مشاهده آنها به صورت روشی برای آرشیوسازی داده‌هایی است که این اسناد در بردارند. این امر مستلزم مدل‌های جدیدی برای ذخیره‌سازی و استخراج داده‌هایی است که از مدل‌های متمایل به تحلیل‌های داده‌های داده‌های ساختار یافته، به جای تحلیل‌های داده‌های ساختار نیافته، پیروی می‌کنند.

وب‌های ملی: چه ارزشی در ایجاد یک وب ملی زمانی که وب پدیده‌ای فراتر از مرزهاست، وجود دارد؟ تلاش‌هایی در حوزه آرشیو- در جعبه- برای اجرا در سطوح ملی، با توجه به محدودیت‌های بودجه‌ای و قانونی، اساسی است. در حال حاضر، در بریتانیا، کتابخانه بریتانیا، خودش را برای پیش‌بینی قوانینی تأثیرگذار بر روی آرشیوسازی فضای وب بریتانیا به صورت کتابخانه واسپاری برای همه انتشارات بریتانیا آماده می‌کند. در می ۲۰۱۱، وزارت پژوهش و اختراع دانمارک، تصمیم به ایجاد زیرساخت‌های پژوهشی ملی در حوزه‌های پژوهشی مختلف (علوم طبیعی، علوم انسانی) که متمرکز بر استفاده از ابزارهای تحلیلی منابع وب آرشیوی خواهد بود.

چالش فوری: عدم شفافیت نحوه استفاده پژوهشگران از آرشیوهای ملی، بسیاری از آنها هنوز در مراحل طراحی هستند. ما در مورد یکی از مهم‌ترین مواردی که باید انجام شود، یعنی تعامل دامنه پژوهشگران با تخصص‌هایی نه فقط در پژوهش اینترنت، بلکه در زمینه‌هایی نظیر جامعه‌شناسی، علوم سیاسی، و دیگر حوزه‌های علوم اجتماعی، فیزیک و سایر علوم، هنر و علوم انسانی، بحث خواهیم کرد. همان‌طور که این زیرساخت‌ها به منظور انعکاس نیازهای پژوهشگران ملی در مجموعه‌های ایجاد شده، طراحی می‌شوند. این فرآیندی زمان‌بر است و تعامل خبرگان هر حوزه می‌تواند مشکل باشد. با وجود

تصویر 23- مشاهده گر n-gram کتاب های گوگل

چالش: برای بخش‌هایی از اینترنت که حاوی داده‌ها به جای اسناد هستند، تغییر نگرش از این که آرشیوهای وب برای نگهداری اسناد هستند به سمت مشاهده آن‌ها به صورت روشی برای آرشیوسازی داده‌هایی است که این اسناد در بردارند. این امر مستلزم مدل‌های جدیدی برای ذخیره‌سازی و استخراج داده‌هایی است که از مدل‌های متمایل به تحلیل‌های داده‌های داده‌های ساختار یافته، به جای

وب های ملی: چه ارزشی در ایجاد یک وب ملی زمانی که وب پدیده ای فراتر از مرزهاست وجود دارد؟ تلاش هایی در حوزه آرشیو - در جعبه-برای اجرا در سطوح، ملی با توجه به محدودیت های بودجه ای و، قانونی اساسی است در حال حاضر، در بریتانیا، کتابخانه بریتانیا خودش را برای پیش بینی قوانینی تأثیر گذار بر روی آرشیوسازی فضای وب بریتانیا به صورت کتابخانه واسپاری برای همه انتشارات بریتانیا آماده می. کند در می، 2011 وزارت پژوهش و اختراع دانمارک، تصمیم به ایجاد زیر ساخت های پژوهشی ملی در حوزه های پژوهشی مختلف (علوم طبیعی علوم انسانی) که متمرکز بر استفاده از ابزارهای تحلیلی منابع وب آرشیوی خواهد بود.

چالش فوری: عدم شفافیت نحوه استفاده پژوهش گران از آرشیو های ملی. بسیاری از آن ها هنوز در مراحل طراحی هستند ما در مورد یکی از مهم ترین مواردی که باید انجام شود، یعنی تعامل دامنه پژوهشگران با تخصص هایی نه فقط در پژوهش، اینترنت بلکه در زمینه هایی نظیر جامعه شناسی، علوم سیاسی و دیگر حوزه های علوم اجتماعی فیزیک و سایر علوم هنر و علوم انسانی، بحث خواهیم کرد. همان طور که این زیرساخت ها به منظور انعکاس نیازهای پژوهشگران ملی در مجموعه های ایجاد شده، طراحی می شوند. این فرآیندی زمانبر است و تعامل خبرگان هر حوزه می تواند مشکل باشد. با وجود

این شکست اجرای آن احتمال کسب کاربرد گسترده زیر ساخت های جدید را کاهش می دهد.

نتایج: مسیر پیش رو

آن چه گفته شد، فقط تعداد معدودی از مواردی است که گروه کوچک ما توانست به آن فکر کند و موارد بسیاری هنوز وجود دارد. بعضی از مباحث عمومی است زیرا فنون گام به گام ویژه ای برای مشخص نمودن مجموعه وب یا روندهای تحلیلی از اینکه چگونه محتوای وب در طول زمان تغییر می کند نیاز به منابعی از پروژه تحقیقاتی خاص، یک تیم از متخصصان حوزه و مجموعه های مرتبط برای آزمون این روش ها خواهد داشت، بنابراین ما گوی جادویی نخواهیم داشت با این حال تعدادی از چالش های عمومی را نشان داده ایم که پژوهشگر علاقه مند به کار در حوزه آرشیوهای وب با آن مواجه است یکی از موارد اصلی که در طول زمان حاصل می شود و دوباره در مصاحبه ها و بحث ها وجود دارد، فقدان پایداری و رابطه ای کاربر پسند برای ایجاد آرشیوهای وب و یکبار ساخته شده به منظور دسترسی و تحلیل داده های موجود در آن هاست.

در دراز مدت، ما امیدوار به نتایج دیگری هستیم که ممکن است از این تلاش ها حاصل شود. برای مثال، می توانیم گروه کاری پست هاگ (1) را تصور کنیم که برای توسعه کارگاه بحث و ایده های مورد استفاده احتمالی در آینده و تمرکز بر توسعه ابزاری ایجاد شد. این گروه نه تنها برای اعضای IIPC، بلکه برای انواع جوامع پژوهش گرانی که با جامعه آرشیوی وب در تعامل نبوده و بیش تر تمایل به استفاده از وب آرشیوی را دارند آگاهی رسانی خواهد کرد. مثال ها شامل پژوهش گران اینترنت (مانند اعضای 2) AOIR) دانشمندان اطلاعات (مانند 3) IFIP و 4) ASIS T) و طیف فهرست ها و انجمن های علاقه مند به علوم انسانی رقومی هستند. ما قویاً توصیه می کنیم که IPC نمایندگانی را به نشست های سالانه سازمان هایی از این نوع و پانل ها رهنمون می سازد و کارگاه هایی را برای تعامل پژوهشگران با امکانات آرشیوهای وب سازماندهی کند. ما چند ایده را به صورت برجسته مطرح ساخته ایم اما آن جوامع موارد بیش تری را رانند آماده نمایند. برای آمدن آن ها به IIPC منتظر نمایید. IPC باید به سمت آن ها برود.

ایده دیگر برای آینده فعالیت، نوعی برنامه مشترک کدنویسی نرم افزاری (5) است که برنامه نویسان رایانه ای هکرها با همدیگر و با پژوهشگران برای 2 الی 3 روز گردهم آمده و به ایجاد دسترسی به داده های آرشیوی وب توجه و همکاری می کنند آن ها باید بتوانند به گروه ها و وظایف مرتبط با کشف رهیافت های خلاقانه و ابتکاری توجه کنند تا بتوانند داده ها و ابزارهای موجود را برای ایجاد سریع ابزارها و رابطه ای جدید، به کار گیرند پژوهش گران انتخاب خواهند شد به دلیل این که آن ها سؤالاتی دارند

ص: 144

1- (گروهی کاری در لاهه) post-Hague

2- AoIR)

3- IFIP

4- ISAST

5- Hackathon: برنامه ای است که طی آن تعدادی از برنامه نویسان به صورت مشترک و فشرده در یک دوره کوتاه زمانی اقدام به کدنویسی برای برنامه خاص می نمایند منبع قابل دسترسی در:

<http://www.techopedia.com/definition/23193/hackathon>

که تمایل به پاسخ‌گویی به آن‌ها دارند و برنامه نویسان رایانه‌ای همه مهارت‌شان را برای کمک به آن‌ها در رسیدن (نیل یا نزدیک شدن) به اهداف پژوهشی‌شان به کار می‌برند دوباره برنامه نویسان رایانه‌ای که با داده‌های وب پویا کار می‌کنند، مهارت‌هایی برای اجرای طرح‌های خلاقانه با ابزارها دارند که نتایج بسیار بیشتری از انتظار برای درخواست از آرشیوهای وب به سمت آن‌ها گسیل خواهد کرد.

آیا ما به نیروانا خواهیم رسید یا محکوم به آخر الزمان خواهیم شد، شد یا با منحصر به فردی (فنی) جایگزین خواهیم شد یا اینکه با آرشیوهای غبار آلوده مواجه خواهیم شد؟ هیچ راهی برای دانستن آن نداریم با این حال ما در نقطه‌ای هستیم که این سؤال را ایجاب می‌کند امروز برای اطمینان از این که در آینده به چه چیزی می‌توانیم دسترسی داشته باشیم چه گام‌هایی می‌توانیم برداریم؟ تجمیع تصمیمات به ندرت قابل توجه، ساده، نبود اما بخشی از تلاشی است برای تضمین این که آرشیوهای وب قوی پایدار قابل دسترسی، ارزشمند و بالاتر از همه توسط پژوهش‌گران آینده قابل استفاده باشند.

منابع

Conover, M. D., Ratkiewicz, J., Francisco, M., Gonçaves, B., Flammini, A., Menczer, F. (2011). .1 Political Polarization on Twitter. Paper presented at the ICWSM: International Conference on Weblogs and Social Media 2011, Barcelona

Conover, M. D., Ratkiewicz, J., Gonçaves, B., Flammini, A., Menczer, F. (2011). The Echo Chamber. .2 Paper presented at the Journal of Information Technology Politics Conference 2011: The Future of Computational Social Science, Seattle

Dougherty, M., Meyer, E. T., Madsen, C., Van den Heuvel, C., Thomas, A., Wyatt, S. (2010). Researcher .3 Engagement with Web Archives: State of the Art. Report. London: JISC. Retrieved from <http://ssrn.com/abstract=1714997> and <http://ie-repository.jisc.ac.uk/544>

Garfinkel, S. Cox, D. (2009, 9–11 February). Finding and Archiving the Internet Footprint. Paper .4 presented at the First Digital Lives Research Conference: Personal Digital Archives for the 21st Century, London

.5 (Gazan, R. (2008). Social annotations in digital library collections. D-Lib Magazine, 14(11/12)

Hindman, M. (2007). "Open-source politics" Reconsidered: Emerging Patterns in Online Political .6 Participation. In V. Mayer-Schönberger D. Lazer (Eds.), Governance and information technology: From electronic government to information government (pp. 183–207). Cambridge: The MIT Press

- Hogan, B. (2010). Analyzing Facebook Networks. In D. Hansen, M. Smith B. Schneiderman (Eds.), . 7
Analyzing Social Media Networks with NodeXL. New York, NY: Morgan Kaufman
- Jasra, M. (2011, 3 February). Reddit Surpasses 1 Billion Monthly Page Views Retrieved 30 April, 2011, . 8
from <http://www.webanalyticsworld.net/2011/02/reddit-surpasses-1-billion-monthly-page.html>.
(Archived by WebCite® at <http://www.webcitation.org/5yKdMBKNC>)
- Kay, A. (1995). The Best Way to Predict the Future is to Invent it. *Mathematical Social Sciences*, 30, . 9
.326–326
- Klett, M., Manchester, E., Verburg, J. (1984). *Second View: The Rephotographic Survey Project*. . 10
.Albuquerque: University of New Mexico Press
- Kling, R., McKim, G., King, A. (2003). A Bit More to IT: Scholarly Communication Forums as Socio- . 11
Technical Interaction Networks. *Journal of the American Society for Information Science and Technology*,
.54(1), 46–67
- .Kurzweil, R. (2005). *The Singularity is Near: When Humans Transcend Biology*. New York: Viking . 12
- Mayer-SchÖnberger, V. (2009). *Delete: the virtue of forgetting in the digital age*. Princeton, NJ: . 13
.Princeton Univ Press
- Meyer, E. T. (2006). Socio-technical Interaction Networks: A discussion of the strengths, weaknesses . 14
and future of Kling's STIN model. In J. Berleur, M. I. Numinem J. Impagliazzo (Eds.), *IFIP International
Federation for Information Processing, Volume 223, Social Informatics: An Information Society for All?* In
.Remembrance of Rob Kling (pp. 37–48). Boston: Springer
- Meyer, E. T. (2011). *Splashes and Ripples: Synthesizing the Evidence on the Impact of Digital* . 15
.Resources. Report. London: JISC. Retrieved from <http://ssrn.com/abstract=1846535>
- Meyer, E. T., Eccles, K., Thelwall, M., Madsen, C. (2009). *Final Report to JISC on the Usage and* . 16
*Impact Study of JISC-funded Phase 1 Digitisation Projects the Toolkit for the Impact of Digitised Scholarly
Resources (TIDSR)*. Retrieved from [http://microsites.oi.ox.ac.uk/tidsr/system/files/TIDSR-Final-Report-
20July2009.pdf](http://microsites.oi.ox.ac.uk/tidsr/system/files/TIDSR-Final-Report-20July2009.pdf)
- .Moretti, F. (2005). *Graphs, Maps, Trees: Abstract models for a literary history*. London: Verso Books . 17

- .Moretti, F. (2011). Network Theory, Plot Analysis. *New Left Review*, 68, 80–102 .18
- Olston, C., Reed, B., Srivastava, U., Kumar, R., Tomkins, A. (2008). Pig Latin: A Not- So-Foreign .19
Language for Data Processing. Paper presented at the ACM SIGMOD'08 Conference, Vancouver, BC,
.Canada
- Schroeder, R. (2011). *Being There Together: Social Interaction in Shared Virtual Environments*. New .20
.York, NY: Oxford University Press USA
- Schroeder, R. Meyer, E. T. (2009). An Emerging Global Brain: How the Internet is Revolutionising .21
Scientific Research. *Britain in 2009 (Economic Social Research Council Annual Magazine)*, 113.27
- Tanner, S. (2010). *Inspiring Research, Inspiring Scholarship*. Report. London: JISC. Retrieved from .22
[http://www.jisc.ac.uk/media/documents/programmes/digitisation/12page
.efinaldocumentbenefitssynthesis.pdf](http://www.jisc.ac.uk/media/documents/programmes/digitisation/12pagefinaldocumentbenefitssynthesis.pdf)
- Tanner, S. Deegan, M. (2011). *Inspiring Research, Inspiring Scholarship: The value and benefits of .23
digitised resources for learning, teaching, research and enjoyment*. Report. London: JISC. Retrieved from
[http://www.kdcs.kcl.ac.uk/fileadmin/documents/ Inspiring-Research-Inspiring-Scholarship-2011-Simon
.Tanner.pdf](http://www.kdcs.kcl.ac.uk/fileadmin/documents/Inspiring-Research-Inspiring-Scholarship-2011-SimonTanner.pdf)
- Thomas, A., Meyer, E. T., Dougherty, M., Van den Heuvel, C., Madsen, C., Wyatt, S. (2010). .24
Researcher Engagement with Web Archives: Challenges and Opportunities for Investment. Report. London:
.JISC. Retrieved from <http://ssrn.com/abstract=1715000> and <http://ie-repository.jisc.ac.uk/543>
- van den Heuvel, C. (2009). MAPS: Manuscript Map Annotation and Presentation System: Linking .25
formal ontologies with social tagging to (re-) construct relationships between manuscript maps and
contextual documents. *Digital Humanities 2009 (University of Maryland, Maryland Institute for
.Technology in the Humanities (MITH) Abstracts)*, 138– 141
- Von Ahn, L., Maurer, B., McMillen, C., Abraham, D., Blum, M. (2008). reCAPTCHA: Human-Based .26
.Character Recognition via Web Security Measures. *Science*, 321(5895), 1465–1468
- Williams, D., Yee, N., Caplan, S. E. (2008). Who plays, how much, and why? Debunking the .27
stereotypical gamer profile. *Journal of Computer-Mediated Communication*, 13(4), 993–1018. doi:
10.1111/j.1083-6101.2008.00428.x

فصل دوم: تجارب جهانی و مسائل بومی در آرشيو سازی وب

اشاره

ص: 149

کتابخانه ملی، استرالیا مؤسسه راهبر آرشیو و حفاظت رقومی در استرالیاست. آرشیو پاندورا (1)، که بیش از 10 سال مخزن آرشیو منابع وب استرالیا بوده است، سیستمی کامل و در حال پیشرفت و توسعه مداوم است. پانداس PANDAS، سیستم مدیریت، آرشیوی که آرشیو را پشتیبانی می کند از سال 2007، در حال سومین بازبینی خود می باشد فعالیت های دیگر آرشیو وب از جمله گردآوری سالانه دامنه استرالیا و استفاده از آرشیو - آی تی (2) با همکاری آرشیو اینترنت اداره می شوند.

این مقاله به بررسی وضعیت کنونی آرشیو وب در استرالیا می پردازد و نشان می دهد که چگونه کتابخانه ها خدمات خود را با ورود فزاینده مواد برخط در مجموعه های شان وفق می دهند سالیان متمادی تصور می شد که با آرشیو کردن فقط می توانیم نمونه ای کوچک اما روایت گر از اینترنت را ضبط کنیم. امروزه شکاف بین آن چه موجود است و آن چه می تواند آرشیو شود در حال کاهش است. درعین حال هر چه آرشیو ها بیشتر می شوند و توانایی ما در آرشیو کردن افزایش می یابد بیش تر با فناوری های جدید و برنامه های کاربردی وب 2/0 درگیر خواهیم بود به عنوان مثال در انتخابات فدرال سال 2007 که تعداد وسیعی از سایت های متعامل نظیر MySpace Kevin و YouTube آرا را آرشیو کردند، معلوم شد که آرشیوکننده های استرالیایی وب به سازگاری ادامه می دهند و با چالش های جدید روبه رو می شوند.

ص: 150

*آرشیو وب در دنیای وب 2/0 شعبه آرشیو وب و حفاظت رقومی کتابخانه ملی استرالیا (1)

ادگار کروک (2) | ترجمه: مرجان هادی زاده (3)

مقدمه

مقاله حاضر آرشیو وب در استرالیا را مورد بحث و بررسی قرار می دهد. از آن جا که کتابخانه ملی استرالیا نقش اصلی را در این زمینه دارد و همچنین با توجه به این که نویسنده خود در آن جا شاغل است این مقاله می تواند تأثیر زیادی بر مسائل این مؤسسه داشته باشد. لازم به ذکر است که در استرالیا پروژه های آرشیو وب دیگری نیز وجود دارند نظیر پروژه Our Digital Island در منطقه تاسمانی (<http://odi.statelibrary.tas.gov.au/>) و پروژه خدمات Territory Stories (<http://www.territorystories.nt.gov.au/>) در منطقه شمالی.

اخیراً کتابخانه ملی، استرالیا برای اجرای آرشیو وب از سه روش شناسی مختلف استفاده می کند. آرشیو انتخابی در پاندورا (آرشیو وب استرالیا)، طی قراردادی با آرشیو، اینترنت گردآوری کل دامنه و بهره برداری از خدمات آرشیو - آی تی را با همکاری یکدیگر انجام داده اند. با این روش به سمت ایجاد مجموعه کامل و جامع نشریات برخط استرالیا حرکت می کنیم البته با افزایش چالش های نوین فناوری کتابخانه مجبور است برای تداوم این امر مهم روش هایی برای تعدیل طرح های آینده اتخاذ کند؛ نظیر محدوده مجموعه ها و گسترش همکاری های جدید

ص: 151

Web Archiving in a WEB 2/0 World -1

Edgar Gruk -2

3- کارشناس ارشد سازمان اسناد و کتابخانه ملی ایران

آرشیو پاندورا (<http://pandora.nla.gov.au>)، انتشارات وب استرالیا را از سال 1996 بایگانی کرده است؛ و در همان زمان به عنوان مخزن مورد قبول در میان تمام ملت ها معرفی شد. اصلی ترین موفقیت پاندورا ساخت شبکه ای مشتمل بر 9 گروه آرشیو استرالیایی بوده است مانند تمام کتابخانه های دولتی قاره استرالیا و AIATSIS (مؤسسه استرالیایی مطالعات بومیان و جزیره نشینان تنگه تارس)، آرشیو ملی فیلم و صدا و یادمان جنگ استرالیا، بنیانگذاری پانداس - سیستم گردش کار آرشیو که سیستمی است برای شناسایی مداوم مطالب آرشیو شده و ارتباط با طیف وسیعی از سازمان های نمایه سازی و چکیده نویسی - و همچنین تعداد بی شماری از ناشران استرالیایی؛ به گونه ای که در اول جولای 2008، آرشیو، در برگیرنده 19,307,307 عنوان شامل 531,140,080 فایل بالغ بر 2/2 ترابایت داده بوده است.

آرشیو پاندورا موضوع های زیر را آرشیو کرده است:

- انتشارات برخط منتخب استرالیا در سطح جهانی نظیر مجله های الکترونیکی، انتشارات دولتی و وب گاه های مهم پژوهشی و فرهنگی؛
- خط مشی ها، شیوه ها و دستورالعمل های منتخب برای مجموعه و تأمین دسترسی بلند مدت به مواد موجود در آرشیو؛
- رویکرد ملی مشترک برای آرشیو و حفاظت بلند مدت انتشارات پیوسته استرالیا، شامل مشارکت کتابخانه های دولتی و دیگر مؤسسه های فرهنگی؛
- سیستم آرشیو رقومی (پانداس) برای جمع آوری و بارگذاری ساده و مؤثر انتشارات در آرشیو، ذخیره اطلاعات در مورد آن ها و مدیریت دسترسی عمومی به آن ها؛
- طرحی برای نام گذاری مستمر تمام اشیای موجود در آرشیو و ارائه خدمات با دقت به آن ها؛
- تنظیم مقررات همکاری با مؤسسه های نمایه سازی و چکیده نویسی در آرشیو پاندورا، به منظور استناد و دسترسی دائمی به مواد بایگانی شده و شناسه های ماندگار برای انتشارات نمایه سازی و چکیده نویسی در آن ها؛ و
- محتوا که در اول جولای 2008 بالغ بر 19,307,307 عنوان مشتمل بر 531,120,080 فایل به میزان 2/2 ترابایت داده می باشد

گستره وسیعی از انتشارات وجود دارد نیمی از آن ها از وب گاه های دولت فدرال و ایالتی آرشیو شده اند و نیمی دیگر بازتابی از تنوع کامل فرهنگ و پژوهش استرالیاست نوع انتشارات درون آرشیو می تواند یک تک مدرک PDF یا یک وبگاه کامل سازمانی شامل هزاران فایل باشد در آرشیو، وینوشت، پادکست، و فیلم نیز وجود دارند.

جمع آوری گزینشی هدایت شده انتشارات، وب بر اساس ارزش بلند مدت فرهنگی و پژوهشی، آن ها بدان معناست که تنها بخش کوچکی از دامنه استرالیا آرشیو شده است کتابخانه، از سال 2005، این مطلب را دریافت و از آن زمان با آرشیو اینترنت باشد در آرشیو، وینوشت، پادکست، و فیلم نیز وجود دارند. (<http://www.archive.org>) قرارداد بست که خزش های هدایت شده سالانه را به منظور جمع آوری هر ماده ممکن از دامنه استرالیا، انجام دهد. این

خزش ها هر سال حدود یک ماه صورت می گیرد و مقدار داده هایی که جمع آوری می کند به اندازه ای است که محتوای پاندورا را کم جلوه می دهند؛ به عنوان مثال در سال 2007 از کل گردآوری دامنه 18 ترابایت داده در یک ماه برداشت شد در حالی که پاندورا در 11 سال 2 ترابایت داده جمع آوری کرده است. انتظار می رود جمع آوری سال 2008 یک بیلیون فایل باشد.

عکس

آرشیو وب در دنیای وب ۱۵۳

خزش ها، هر سال حدود یک ماه صورت می گیرد و مقدار داده هایی که جمع آوری می کند به اندازه ای است که محتوای پاندورا را کم جلوه می دهند؛ به عنوان مثال، در سال ۲۰۰۷، از کل گردآوری دامنه ۱۸ ترابایت داده در یک ماه برداشت شد، در حالیکه، پاندورا در ۱۱ سال ۲ ترابایت داده جمع آوری کرده است. انتظار می رود جمع آوری سال ۲۰۰۸ یک بیلیون فایل باشد.

گردآوری دامنه وب استرالیا: تجزیه و تحلیل کمی مقدماتی آرشیو داده، NLA، ۲۰۰۸ -

تاریخ گردآوری دامنه	۲۰۰۵	۲۰۰۶	۲۰۰۷
مدارک (فایل های) تکی خزش شده	۱۸۵۵۴۹۶۶۲	۵۹۶۲۳۸۹۹۰	۵۱۶۰۶۴۸۲۰
کل مدارک (فایل های) خزش شده	۱۸۹۸۲۴۱۱۹	۶۲۱۶۶۴۸۷۶	۵۲۳۵۱۰۹۴۵
میزبان ها	۸۱۱۵۲۳	۱۲۶۰۵۵۳	۱۲۴۷۶۱۴
اندازه داده های خام	۶,۶۹ ترابایت	۱۹,۰۴ ترابایت	۱۸,۴۷ ترابایت
اندازه فایل های ARC فشرده شده	۴,۵۲ ترابایت	۱۰,۴۸ ترابایت	۱۰,۱۸ ترابایت

HTTP://PANDORA.NLA.GOV.AU/DOCUMENTS/AUSCRAWLS.PDFKOERBIN, P.

گردآوری هایی که با استفاده از HERITRIX هدایت شده اند (<http://crawler.archive.org>) بزرگ است، ولی کاملاً جامع نیستند، زیرا فقط یک ماه در هر سال گردآوری می شوند (و مقادیر زیادی می توانند در این فاصله زمانی در اینترنت بیابند و بروند)، آنها از قوانین robots.txt پیروی می کنند، و اگرچه HERITRIX عملکردی قوی دارد، و نگاه هایی وجود دارد که از لحاظ فنی خزش به آنها مشکل است. نقص ها هر چه باشند، گردآوری در اندازه قابل توجهی انجام می گیرد و بنابراین مقادیر زیاد داده گردآوری شده باعث می شود هرگونه تلاشی در جهت ارزشیابی کیفی و نگاه های شخصی مشکل گردد. بنابراین، به طور معکوس، پاندورا، جایی است که می توان مشکلات را درون هر عنوان شناسایی کرد و تحلیل کیفی را انجام داد - در صورت لزوم صفحه ها را گردآوری کرد و تثبیت کرد - که در اینجا ممکن نیست. یکی دیگر از اشکال ها این است که برخلاف پاندورا - که در آن اجازه ناشر برای آرشیو کسب می شود - در اینجا چنین امکانی وجود ندارد. بنابراین، با توجه به قانون کپی رایت (حق تألیف) فعلی استرالیا که طبق آن نشریات برخط شامل واسپاری قانونی نیستند، در حال حاضر، قادر نیستیم آنچه را که حفظ شده به عموم نمایش دهیم. البته، این عبارت به این معنی نیست که هیچ کاری در این آرشیو انجام نمی شود، چرا که دانشگاهیان در حال کار بر روی این داده ها هستند و به نظر می رسد تحقیقات آنها

گردآوری دامنه وب استرالیا: تجزیه و تحلیل کمی مقدماتی آرشیو داده، NLA، 2008 -

گردآوری هایی که با استفاده از HERITRIX هدایت شده اند (<http://crawler.archive.org>) بزرگ است، ولی کاملاً جامع نیستند زیرا فقط یک ماه در هر سال گردآوری می شوند (و مقادیر زیادی را می توانند در این فاصله زمانی در اینترنت بیابند و بروند آن ها از قوانین robots.txt پیروی می کنند، و اگر چه HERITRIX عملکردی قوی دارد وب گاه هایی وجود دارد که از لحاظ فنی خزش به آن ها مشکل است. نقص ها هر چه باشند گردآوری در اندازه قابل توجهی انجام می گیرد و بنابراین مقادیر زیاد داده گردآوری شده باعث می شود هر گونه تلاشی در جهت ارزشیابی کیفی وب گاه های شخصی مشکل گردد.، بنابراین به طور معکوس، پاندورا جایی است که می توان مشکلات را درون هر عنوان شناسایی کرد و تحلیل کیفی را انجام داد- در صورت لزوم صفحه ها را گردآوری کرد و تثبیت کرد - که در این جا ممکن نیست یکی دیگر از اشکال ها این است که بر خلاف پاندورا- که در آن اجازه ناشر برای آرشیو کسب می شود- در اینجا چنین امکانی وجود ندارد، بنابراین با توجه به قانون کپی رایت (حق تألیف) فعلی استرالیا که طبق آن نشریات برخط شامل واسپاری قانونی نیستند در حال حاضر، قادر نیستیم آن چه را که حفظ شده به عموم نمایش دهیم البته این عبارت به این معنی نیست که هیچ کاری در این آرشیو انجام نمی شود، چرا که دانشگاهیان در حال کار بر روی این داده ها هستند و به نظر می رسد تحقیقات آن ها

ارزشمند باشد آرشیو-آی تی سرویس آرشیو وب های میزبانی شده توسط آرشیو اینترنت است. اولین (و تا کنون تنها) سازمانی که در استرالیا از این سرویس استفاده کرده و می کند بخش مجموعه های آسیایی کتابخانه ملی (<http://www.nla.gov.au/asian/asianwebarchive.html>) است.

از آرشیو-آی تی برای جمع آوری مجموعه ای از وب گاه های خارج از کشور استفاده می شود که رویدادهای خاص اجتماعی و سیاسی را ضبط می کنند؛ زیرا انتظار نمی رود که هیچ سازمان دیگری در منطقه این نقش را به انجام برساند گزینه میزبانی نیز از آن رو انتخاب شد که به نظر می رسد می تواند راهی سریع و آسان برای جمع آوری و نگهداری مجموعه ها باشد و به مهارت های فنی نیاز ندارد و میزان زیادی از وقت کارکنان را نمی گیرد. در حالی که خیلی زود معلوم شد فقط قسمتی از این واقعیت دارد چرا که برای ساخت موفقیت آمیز مجموعه ها زمانی بسیار زیادتر از آن چه تصور می شد صرف شد گزینش وب گاه ها برای ، خزش فعالیتی است که اغلب با سو تفاهم همراه است و می تواند به طور شگفت آوری زمان زیادی بطلبد.

در حالی که استفاده از آرشیو آی تی دارای این مزیت است که دیگر در مورد میزبانی و حفظ محتوای جمع آوری شده نگرانی وجود نخواهد داشت. البته برخی مشکلات عمده در مورد عدم کنترل مواد آرشیوی و نمایش توابع وجود دارد. آرشیو آی تی اجازه می دهد تا برخی URL های هسته انتخابی جمع آوری شوند؛ با وجود این، اگر فایل های جمع آوری شده از بین بروند چه جمع آوری بشوند چه جمع آوری نشوند، هیچ راهی برای اصلاح خرابی یا از دست دادن محتوا وجود ندارد که بتوانید با سیستم خود انجام دهید به همین ترتیب، کنترل واقعی یا مالکیت فرآیند نمایش نیز وجود ندارد؛ به طوری که به عنوان مثال یک پیوند به URL هسته که جمع آوری نشده ، هم چنان درون یک مجموعه ظاهر می شود یکی دیگر از اشکال ها این است که اگر شما اشتراک سالیان خدمات خود را قطع کنید مجموعه های شما به مخزن محتوای آرشیو اینترنت عمومی باز می گردد. با وجود این مسائل کتابخانه در نظر دارد به آرشیو با استفاده از این روش ادامه دهد.

سایت های آرشیو شده با استفاده از آرشیو آی تی برای مجموعه های آسیا - کتابخانه ملی استرالیا

Papua New Guinea Government and Research Websites

(<http://archive-it.org/collections/1039>)

وب گاه های انتخابی مؤسسه های دولتی و پژوهشی مهم پاپوا گینه نو که از سال 2008 آرشیو شده اند. برخی سایت های آرشیوی که در زمان ضبط کردن جاری نبودند.

(<http://archive-it.org/collections/918>)

(<http://archive-it.org/collections/1040>)

وب گاه انتخابی بین المللی بین الدولی مرتبط با انتخاب عمومی سال 2007 تایلد

وب گاه های منتخب دولتی، احزاب سیاسی و رسانه های مرتبط با انتخابات ملی کامبوج سال 2008

(<http://archive-it.org/collections/937>)

(<http://archive-it.org/collections/1054>)

وب گاه های منتخب بین المللی مرتبط با شورش راهبان برمه سپتامبر - اکتبر 2007

وب گاه های انتخابی دولتی و غیر دولتی جمهوری دموکراتیک خلق لائوس که از سال 2008 آرشیو شده اند.

ص: 154

وب گاه های انتخابی بین المللی بین دولتی و پاپوآگینه نو مرتبط با انتخابات پارلمانی پاپوآگینه نو در سال 2007. شامل تصاویری از وبگاه کمیسیون انتخابات پاپوآگینه نو، دادخواست انتخابات پاپوآگینه، نو گزارشی از گروه ارزشیابی انتخابات انجمن جزایر مشترک المنافع اقیانوس آرام و بیانیه رسانه ای شفافیت بین المللی (پاپوآگینه نو).

وب گاه های منتخب بین المللی بین دولتی مرتبط با انتخابات ریاست جمهوری و پارلمانی تیمور شرقی در سال 2007 شامل آموزش رأی دهندگان و مواد تبلیغاتی سیاسی

وب گاه های منتخب اندونزیایی

جمع آوری فایل ها

وقتی پاندورا پا به عرصه وجود گذاشت به علت ناتوانی اولیه نرم افزارهای جمع آوری وب، بسیاری از وب گاه ها نمی توانستند جمع آوری شوند وب گاه های با پایه HTML می توانستند جمع آوری شوند، ولی سایت هایی با چنین قالب های ساده ای در ابتدا بسیار مشکل به وجود آوردند تا یک مشکل حل می گردید مشکل دیگری پیدا می شد بنابراین برای آرشویست های وب پیشرفت شگفت انگیز اینترنت جاوا، اسکریپت اپلت ها، شیوه نامه های آبشاری Shockwave flash، و ده هزار فایل دیگر و انواع قالب ها فقط مشکلاتی پس از دیگری بود در حالی که مشکلات اغلب انواع فایل ها بر طرف شده اند، محتوای چندرسانه ای یک مسئله باقی مانده است پیش از این در مورد فایل های پخش (Real Player) و در حال حاضر، درباره پادکست ها مجبوریم نه با پیچیدگی خود فایل ها که با پیچیدگی سیستم های تحویل آن ها مقابله کنیم در حال حاضر این امر در مورد فیلم ها این مسئله لاینحل مانده است.

مجموعه انتخابات فدرالی سال 2007 بزرگ ترین تلاش تاکنون بوده است کتابخانه ملی مسئول آرشیو تمام منابع ملی مرتبط با انتخابات بود از جمله وب گاه های احزاب، گروه های لابی وب گاه های برخی نامزدها، وب نوشت ها، فیلم ها و وب گاه های رسانه ای کتابخانه های ایالتی مسئول جمع آوری وب گاه های نامزدها، احزاب، و رسانه های محلی در ایالت خود بودند در مجموع بیش از 350 وبگاه توسط کتابخانه ملی و همکارانش آرشیو شده است، بسیاری از این سایت ها چندین بار جمع آوری شدند تا محتوای در حال تغییر را ضبط کنند بزرگ ترین چالش در آرشیو کردن این انتخابات تعداد زیاد فیلم ها بود؛ البته مشکلی در خود آن ها نبود، بلکه مشکل در مکانیزم تحویل و فناوری های تعبیه شده در آن ها به منظور استفاده مفید کاربران بود.

رویکردهای متفاوتی برای آرشیو کردن فیلم ها بسته به ماهیت وب گاه ها به کار گرفته شده است. برای وب گاه های عمومی که در آن فیلم های مجزا در یک صفحه وب موجود است فایل های فیلم را جداگانه با استفاده از ابزارهای متفاوت و رایگان موجود بارگیری کردیم (گردآورندگان وب به طور کلی نمی توانند فیلم ها را به طور خودکار جمع آوری کنند)؛ و سپس با استفاده از مبدل های فایل ها را از flv به چیزی کاربر پسند تر مانند قالب Mpeg تغییر دادیم جایی که تعدادی فیلم به یک صفحه وب پیوند داشت، فیلم ها در قالب اصلی خود شان باقی ماندند و یک پخش کننده flv. در فایل های جمع آوری شده نصب گردید، به طوری که فایل ها می توانند به راحتی وب گاه

زنده ارائه شوند وقتی که وب گاه انتخابات یوتیوب

ص: 155

(<http://nla.gov.au/nla.arc-76644> شامل بیش از 700 فیلم) را جمع آوری کردیم، با مراجعه به مهارت های فنی در بخش فناوری اطلاعات توانستیم برای فیلم ها URL ها را از سایت جاری استخراج کنیم آن ها را بارگیری کنیم و تغییرات ضروری را برای گردآوری صفحات وب انجام دهیم.

هیچ یک از این فرآیند ها سریع نبود و همه به مقدار نسبتاً خوبی از مهارت های فنی نیاز داشتند، از جمله کد گذاری مجدد صفحات آرشیو شده با تغییراتی که باید انجام می دادیم تا فیلم ها درون آرشیو قابل پخش باشند. همچنین به خاطر انتخابات حفاظت از اسناد برخط در دولت پیشین نیز به ما سپرده شد. با پیش بینی این امر طرح هایی (همان گونه که برای انتخابات پیشین داشتیم) ساخته بودیم و تمام وب گاه های وزارتی دولت را آرشیو کرده بودیم؛ در حالی که آن ها درست قبل از تاریخ انتخابات در حالت مستحفظ بودند این کار صورت گرفت زیرا گمان می رفت با یک تغییر احتمالی در دولت حذف کلی محتوا از اینترنت وجود داشته باشد همان گونه که در حوزه های قضایی دیگر اتفاق افتاده بود. از آن جا که انتشارات وب برخی بخش های دولتی به ویژه بخش هایی که دارای اسناد تغییر یافته بودند، از دید عموم حذف شده، بود این دور اندیشی اجر نهاده شد.

دستورالعمل های جمع آوری

با توجه به گردآوری دامنه استرالیا توسط کتابخانه ملی و آرشیو انتخابی پاندورا، می توان گفت که حجم قابل توجهی از نشریات استرالیا و یا وب گاه ها آرشیو شده اند اگر چه نمی توانیم بگوییم تا چه حد این مجموعه جامع و کامل است با این حال می دانیم که شکاف های بزرگی باقی مانده است.

ما، وب گاه های تولیدی توسط استرالیایی های خارج از دامنه استرالیا را به طور جامع آرشیو نمی کنیم (اگر چه امیدواریم آرشیو اینترنت بسیاری از آن ها را به دست خواهد آورد). خارج از نشریات مرسوم، اما تا اندازه ای دارای اهمیت بیشتر، ما هم چنین میزان بسیار زیادی از محتوای خلاقانه ای را که توسط افراد تولید شده و در وب گاه ها، وب نوشت ها، دنیای مجازی و سایت های شبکه های اجتماعی فیلم، عکس و هنر میزبانی شده اند، جمع آوری نکرده ایم.

تلاش هایی برای جمع آوری برخی از این محتواها وجود دارد اما پروژه های هدایت شده کوچکی هستند. یکی از آن ها مجموعه رقص استرالیایی است که به دنبال گردآوری نمونه های از رقص استرالیایی به محض قرار گرفتن بر روی وب گاه های مختلف و سایت های میزبان فیلم می باشد.

اگر چه کتابخانه ملی با Flickr قرارداد بسته است و از MySpace و YouTube اجازه آرشیو کردن دریافت کرده ایم کتابخانه فقط بخش کوچکی از این منابع را آرشیو یا گردآوری کرده است. هم چنین، دیگر منابع برخط هر جایی که فعالیت های استرالیایی وجود دارد نظیر فضاهای مجازی (Second Life غیره) و شبکه های اجتماعی (Facebook, Bebo و غیره) نیز آرشیو نشده اند. نخستین دلیل این که معمولاً محتوا دارای حق نشر و حفظ حریم خصوصی است که آرشیو کردن آن را مجاز نمی کند و یا به دلیل ماهیت منابع که آن را خارج از اینترنت عمومی قرار می دهد.

ما به طور جداگانه به عنوان کتابدار نیز می توانیم تفاوتی نسبت به حفظ میراث برخط خودمان به وجود آوریم؛

البته با پذیرش این مسئولیت که می‌توانیم درون سازمان خود یا سازمان مادر از حفظ و پشتیبانی آن چه در وب گاه سازمان قرار گرفته اطمینان حاصل نمایم این، جنبش به ویژه به طور قابل توجهی در دولت و دانشگاه‌ها در مورد حرکت نشریات از چاپی به برخط ادامه می‌یابد آن چه ما پیدا کرده ایم این است که نباید مطمئن باشیم که نشریات برخط در وبگاه ناشران در بلند مدت (و یا حتی کوتاه مدت) در دسترس باقی بمانند دانشگاه‌ها در حال حاضر، مجبور به ایجاد مخازن رقومی برای خروجی فکری خود هستند و در این راه نشریات آن‌ها دسترس پذیر باقی می‌مانند. در حالی که ممکن است انتظار داشته باشیم که وب گاه‌های دولتی دسترسی عمومی خود را حفظ نمایند؛ تجربه نشان می‌دهد که این موضوع همیشه صادق نیست بنابراین اگر نشریه‌ای برای مجموعه و یا کاربران شما ارزشمند است، عاقلانه این است که تلاشی در جهت حفظ دسترسی بلند مدت آن انجام شود.

دستور عمل‌های آینده

تا زمانی که چشم انداز در حال رشد اینترنت وجود دارد فناوری‌های جدید و شکاف‌هایی که تازه شناسایی شده‌اند در گردآوری ما وجود دارند که پرداختن به آن‌ها ضروری است بنابراین، به نظر نمی‌رسد آرشیو وب به منطقه‌ای تبدیل شود که در آن شیوه‌ها و یا پروتکل‌های توسعه مجموعه کاملاً تأسیس و یا برقرار گردد. ما همیشه به شناسایی و جمع‌آوری مطالب نیاز داریم نه این که منتظر بمانیم تا آن‌ها به سوی ما بیایند. وب بسیار پویاست و فناوری بسیار متغیر تعداد ناشران آن قدر وسیع است که بعید به نظر می‌رسد ما قادر به ایجاد سبکی مقیاس پذیر - نظیر سبک سیستم سپرده فیزیکی که نیاکان ما برای مواد چاپی انجام داده‌اند باشیم. هنگامی که کتابخانه ملی برای اولین بار شروع به آرشیو کردن وب کرد، ابزار و نهادهای بسیار کمی در خارج وجود داشت که ما می‌توانستیم برای یادگیری با آن‌ها کار کنیم. در نتیجه، باید مکانیسم‌ها و ابزارها را خودمان اختراع می‌کردیم برای این، منظور کتابخانه ملی در درون خودش سیستم مدیریت آرشیوی را - PANDAS - ساخت این سیستم در حال حاضر، در سومین و طبق تصور ما آخرین مرحله خود است چرا که کتابخانه دیگر نمی‌تواند به طور مستقل چنین سرمایه‌گذاری را برای توسعه متحمل شود. با این حال در حال حاضر که آرشیو وب در طیف وسیعی از نهاد‌های بین‌المللی به طور عملی ایجاد شده، است شرکای زیادی وجود دارند که می‌توان این بار را با آن‌ها تقسیم کرد. کنسرسیوم بین‌المللی حفاظت از اینترنت (1) که اعضای آن متشکل از تمام کتابخانه‌های ملی راهبر آرشیو وب از جمله کتابخانه ملی استرالیا و دیگر نهاد‌های مرتبط می‌باشد این توسعه را هدایت و راهبری می‌کند با این، روش کتابخانه می‌تواند به نقش راهبری خود در آرشیو، وب از طریق تطابق با ابزارهای در حال توسعه و سازگاری با کتابخانه‌ها و مؤسسه‌های همکار ادامه دهد.

شرح حال مختصری از پدید آورنده

ادگار کروک از 1999، در کتابخانه ملی استرالیا کار می‌کند. او از سال 2000 بر روی پاندورا: آرشیو وب استرالیا کار می‌کند. قبل از آن در کند. قبل از آن در کتابخانه عمومی ATC کار می‌کرده است.

ص: 157

همان گونه که در بسیاری از متون نیز مطرح است زبان فارسی به واسطه ویژگی های خاص خود چه از نظر رسم الخط و چه از نظر صرف و معنا با چالش های منحصر به فردی در زمینه ذخیره و بازیابی اطلاعات رو به روست در مقالات مختلف به این مشکلات در قالب ها و سطوح مختلف بدون توجه به نوع شناسی این چالش ها اشاره شده است. به عبارت بهتر کمتر مقاله ای در زبان فارسی در حوزه اطلاع رسانی از دیدی مبتنی بر علم زبان شناسی مدرن به بحث دسته بندی این مشکلات و پاسخ دهی به آن ها پرداخته است. نوشتار حاضر بر آن است تا عمده ترین چالش های مطرح این زمینه را در سه گروه رسم الخط مسائل صرفی و مسائل معنایی مورد اشاره قرار داده و پس از بیان نمونه هایی در پیوند با هر یک از این چالش ها و ارائه ساختواره ای درختی از انواع مسائل مطرح در این سه گروه با نگاهی به پژوهش های صورت گرفته در زمینه زبان فارسی و عربی به ارائه راهکار هایی جهت حل این مسائل پردازد.

کلیدواژه ها: چالش های صرفی زبان فارسی؛ چالش های رسم الخط زبان فارسی؛ چالش های معنایی زبان فارسی؛ بازیابی اطلاعات

شعله ارسطو پور (1) | فاطمه احمدی نسب (2)

درآمدی بر مشخصه‌های زبان فارسی

از دیدگاه دستور زبان زایشی (3)، زبان یا دانش زبانی قوه مستقلی از قوای ذهنی به شمار می‌رود که از دیگر قوای شناختی انسان مستقل بوده و خود نیز دارای بخش‌های مستقل نحو، معنا، واژگان و بخش واجی است به عبارت دیگر این رویکرد یک رویکرد حوزه‌ای (4) به زبان است که بین دانش زبانی و کاربرد آن تمایز قائل است (دبیر مقدم 1383، 18-19). در همین راستا در این نوشتار نیز منظور از زبان فارسی دانش زبانی فارسی‌زبانان ایرانی و هم‌چنین تبلور این دانش زبانی در قالب گونه‌نوشتاری فارسی رایج در ایران یعنی خط فارسی است. زبان فارسی یکی از زبان‌های هندو اروپایی و از شاخه زبان‌های ایرانی جنوب غربی است که زبان رسمی ایران و تاجیکستان و یکی از دو زبان رسمی افغانستان است (کامری 1990، 13-16) این در حالی است که تفاوت‌هایی میان این انشقاق‌ها وجود دارد مثلاً فارسی تاجیکی و فارسی رسمی ایرانی از لحاظ دستوری یکسان هستند

ص: 159

1- دکترای علم اطلاعات و دانش‌شناسی مرکز منطقه‌ای اطلاع‌رسانی علوم و فناوری 2 arastoopoor@ricest.ac.ir

2- دکترای زبان‌شناسی همگانی مرکز منطقه‌ای اطلاع‌رسانی علوم و فناوری ahmadinasabricest.ac.ir

3- Generative grammar

4- Modular

و فقط مقداری تفاوت های واژگانی دارند، اما از لحاظ نظام نوشتاری کاملاً متفاوت هستند چرا که خط فارسی ایرانی از الفبای عربی اقتباس شده است و خط فارسی تاجیکی الفبای سیریلیکی دارد. از دیدگاه رده شناسی (1)، زبان فارسی از برخی ویژگی های سه رده زبانی زبان های تحلیلی (2)، زبان های تصریفی (3) و زبان های پیوندی (4) برخوردار است. این تمایز سه گانه بر اساس ساختار واژگانی زبان تعریف می شود و از آن جا که ذخیره و بازبازی اطلاعات و نمایه سازی غالباً بر اساس واژه صورت می گیرد به نظر می رسد که توجه به این تمایز از اهمیت ویژه ای برخوردار باشد. زبان تحلیلی زبانی است که واژه های آن بر اساس نقش (فاعلی، مفعولی و...) هیچ گونه تغییری نکند و صرفاً جایگاه آن در جمله نشان دهنده نقش واژه باشد کریستال (2008، 256) در زبان های، تصریفی واژه غالباً از چند تکواژ (5) تشکیل می شود که نشان دهنده روابط دستوری است در این زبان ها جایگاه واژه در جمله نشان دهنده نقش واژه نیست بلکه تصریف و حالت (6) واژه نقش واژه را به تصویر می کشد (همان، 244) و بالاخره زبان پیوندی زبانی است که واژه های آن غالباً از بیش از یک تکواژ تشکیل شده اند و هر کدام از تکواژها نشان دهنده یک نقش یا رابطه دستوری است (همان، 16). البته لازم به ذکر است که تمایز زبان تصریفی و پیوندی در آن است که تقطیع و تعیین تکواژهای واژه در آن امکان پذیر است و به راحتی می توان مرز آن ها را تعیین نمود اما در زبان های تصریفی واژه غالباً صرف شده و تکواژها نه به صورت منظم و به ترتیب بلکه به شکل تجمعی در یک صورت کلمه (7) ظاهر می شوند صورت های مفرد و جمع واژگان روایت / روایات عاقبت / عواقب مدرسه / مدارس و ده ها مثال دیگر که البته وام واژه هایی (8) از زبان عربی هستند نشان دهنده وجود تصریف در زبان فارسی است. است از دیگر سو صد ها مثال مانند واژه دانشگاهی (دان + ش + گاه + ی) نشان گر وجود ویژگی زبان پیوندی در فارسی است اگر چه تحقیق جامعی در تمایل بیش تر زبان فارسی به هر یک از این سه رده صورت نگرفته اما به نظر می رسد که زبان فارسی به زبان تحلیلی نزدیک تر بوده و از ترتیب واژه برای نمایش نقش های دستوری بیش از شیوه های دیگر بهره می برد. از بحث های مرتبط با بحث توالی اجزای واژگانی می توان به بحث هسته (9) واژه مرکب اشاره نمود در زبان فارسی واژه مرکب از نظر معنایی و بر اساس هسته واژگانی خود به چهار دسته برون، مرکز درون مرکز (هسته آغازی، هسته پایانی) دو سویه و متوازن تقسیم می شود. در واژه مرکب برون مرکز هیچ کدام از اجزای واژه هسته را تشکیل نمی دهند برای مثال «آب سیاه» نام یک بیماری است و نه آب است و نه سیاه رنگ در واژه مرکب درون مرکز یکی از اجزای واژه هسته است و به دو نوع هسته - آغازی و هسته پایانی

ص: 160

Typology -1

Isolating/analytic -2

Inflectional -3

Agglutinating -4

Morpheme -5

Case -6

Word form -7

Loan word -8

Head -9

تقسیم می‌شود. برای مثال «موش خرما» هسته آغازی و «مویرگ» هسته پایانی است. در واژه مرکب دوسویه هر دو جزء به یک مرجع واحد اشاره می‌نماید مانند «سرباز - معلم» و بالاخره اینکه در واژه مرکب متوازن هر دو جزء به یک اندازه در ساخت کلمه مرکب نقش دارند؛ مانند «دخل و خرج» شقاقی (123، 1386-124).

نکات پیش گفته تنها بخشی از مهم‌ترین ویژگی‌های زبان فارسی است که به صورت بالقوه قابلیت ایجاد بزرگ‌ترین چالش‌ها را در بحث ذخیره و بازیابی اطلاعات دارد عمده‌ترین دلیل چنین امری آن است که ماشین قابلیت تشخیص ترتیب و نقش دستوری را به خودی خود ندارد و همین امر منجر به بازیابی منابعی می‌شود که بعضاً مورد نیاز کاربر نبوده و عملاً پاسخی به نیاز وی نخواهند بود. بنابراین بی‌شک بعد از زبان شناختی فرایند ذخیره و بازیابی اطلاعات یکی از مهم‌ترین و در عین حال برانگیزترین مسئله طراحی و کاربرد پذیری پایگاه‌ها است بازیابی اطلاعات با مسائل زبانی گره خورده و نمی‌توان بدون توجه به زبان مبدأ و ویژگی‌های آن نظام‌های بازیابی اطلاعات را در جهت افزایش جامعیت و مانعیت از یک سو و بهبود ربط از دیگر سو ارتقا بخشید. این نوشتار بر آن است تا از سه بعد رسم الخط صرف و معنا به بررسی مشکلات و چالش‌های زبان فارسی پرداخته و به صورت همزمان به نمونه‌هایی عینی از این مشکلات و نتیجه عدم توجه به آن‌ها در سطح وب و برخی از پایگاه‌های مقالات فارسی اشاره کرده و به آسیب‌شناسی زبان و خط فارسی در رابطه با این حوزه بپردازد.

پیشینه پژوهش

تلاش در جهت اصلاح و تقویت زبان و خط فارسی را می‌توان به تأسیس فرهنگستان اول نسبت داد. در این نوشتار به منظور پرهیز از طولانی شدن بحث صرفاً به برخی از تلاش‌های اخیر اشاره شده است. به عنوان نمونه آشوری (1375) برای اصلاح خط فارسی پیشنهاد می‌دهد که صورت‌های صرفی زمان حال فعل بودن یعنی «آم»، «ای»، «است»، «ایم»، «اید»، «اند» و ضمائر متصل «آم»، «ات»، «آش»، «مان»، «تان»، «شان» جدا از کلمه‌های قبل از خود نوشته شوند معصومی همدانی (1381) در مقاله‌ای تحت عنوان خط فارسی و رایانه یکسان‌سازی رسم الخط فارسی را برای استفاده از خطایاب امکان جستجوی واژه در متن و تهیه نمایه ضروری می‌داند وی راه حل را در تهیه یک دستورالعمل واحد برای خط فارسی دانسته و با اشاره به دستورالعمل فرهنگستان می‌نویسد: «متأسفانه میزان آزادی که این دستور خط به استفاده‌کنندگان داده به قدری است که می‌توان گفت کاربرد آن در مواردی بر تشنگی موجود خواهد افزود». اسلامی (1381) اعمال اصلاحاتی را در خط فارسی ضروری می‌داند؛ اصلاحاتی از قبیل وضع یک نشانه اصلی برای نمایش کسره اضافه استفاده از نشانه‌های متفاوت برای نمایشی نکره و ی اسم-ساز و صفت ساز سرهم نویسی کلمات غیر بسیط، قرار دادن علامت‌های جداگانه برای واژه بست (1) های ربطی فعل بودن و ضمائر ملکی، و بالاخره نشان

دادن اسامی خاص محقق زاده و زارعیان (1383) با الگوگیری از زبان های لاتین چپ نویسی کاهش حروف از طریق جدانویسی قائل شدن به دو شکل کوچک و بزرگ حروف را در مورد زبان فارسی پیشنهاد می نمایند. علاوه بر این نگارش حروفی که خوانده شده ولی نوشته نمی شوند، ایجاد علامتی برای کسره اضافه و عدم تمایز بین «آ» و «ا» را نیز در جهت سازگار ساختن زبان فارسی با محیط های رایانه ای ضروری می دانند اسلامی (1386) نقطه دار بودن برخی از نویسه های فارسی را یکی از ایرادهای خط فارسی می داند؛ چرا که علاوه بر دشوار نمودن املا، کلمات این ایراد اساسی را دارد که استفاده از برنامه نشانه خوان نوری (1) را برای متون فارسی دشوار می نماید. صفار مقدم (1386) فاصله گذاری را در ترکیبات فارسی ضروری دانسته و مفصلاً به مبحث فاصله گذاری درون کلمه ای و برون کلمه ای در انواع کلمات فارسی پرداخته است. گل تاجی و بذرگر (1389) با استفاده از یک سیاهه 17 کلید واژه ای مشکلات ریخت شناسی خط فارسی را در سه پایگاه اطلاعاتی مرکز منطقه ای اطلاع رسانی علوم و فناوری پژوهشگاه اطلاعات و مدارک علمی ایران و جهاد دانشگاهی مقایسه کرده و نتیجه می گیرند که هیچ کدام از پایگاه های مذکور به ور جامع به این مسائل توجه نداشته اند.

همان طور که از مرور آثار پیشین پیداست برخی از پژوهش ها به ضرورت اصلاح خط فارسی با توجه به محیط های رایانه ای و برخی به آموزش زبان فارسی توجه داشته اند که برخی از پیشنهادات مطرح شده در این راستا افراطی بوده و در صورت اجرا چهره خط فارسی را کاملاً دگرگون می نمایند برای مثال چپ نویسی استفاده از حروف بزرگ و کوچک و یا نگارش صداهایی که خوانده می شود. خط دارای ماهیتی محافظه کارانه بوده و به راحتی تغییر را نمی پذیرد. علاوه بر این تغییر آن به راحتی امکان پذیر نیست و به بسیج امکانات مادی و معنوی عظیمی نیاز دارد به نظر می رسد که در راستای حفظ گنجینه زبان فارسی باید خط فارسی را به همین صورت حفظ کرد و تنها با اعمال برخی اصلاحات و پیروی از دستور العملی واحد در جهت یکسان سازی خط فارسی از یک سو و کسب توانمندی های بیش تر در حوزه پردازش زبان طبیعی از دیگر سو به سمت رفع مشکلات بازیابی اطلاعات در سطح وب حرکت کرد.

رسم الخط فارسی و بازیابی اطلاعات

در مورد مسائل مرتبط به خط فارسی به طور مطلق و همچنین ارتباط آن با بازیابی اطلاعات آثار متعددی به نگارش درآمده است که هر کدام از آن ها به گوشه ای از این مشکل عمده توجه داشته است در این نوشتار تنها به دو مسئله فاصله گذاری و نگارش «الف» پرداخته می شود یکی از راه حل های ارائه شده برای یکسان سازی خط فارسی و بهبود بازیابی اطلاعات فارسی در وب پیشنهاد عدم تمایز صورت های مختلف الف آغازی است یعنی «آ» و «ا» به یک صورت یعنی «ا» نوشته شود (اسلامی، 1381). اما این راهکار نمی تواند مفید باشد و به ابهام معنایی منجر می شود. برای مثال در پایگاه RICEST بین «آ» و «ا» تمایزی وجود ندارد و باعث می شود که جستجوی کلید واژه های آسم

(بیماری تنفسی) و اسم (مقوله دستوری) در موتور جستجوی جامع به نتایج کاملاً یکسان منجر شود و حال آن که منطقاً این دو کلیدواژه باید به بازیابی نتایج کاملاً متفاوتی منجر شود. به عبارت بهتر، قائل شدن این تمایز سطح ربط نتایج بازیابی شده به میزان چشم‌گیری افزایش می‌یابد. مثلاً در پایگاه مگ ایران این تمایز در نظر گرفته شده است و جستجوی کلیدواژه‌ها به نتایج متفاوت و مرتبط به کلیدواژه مورد نظر می‌انجامد. (1)

چالش عمده دیگر در پیوند با رسم الخط فارسی متوجه نحوه فاصله‌گذاری است. همان‌طور که می‌دانیم اکثر نویسه (2) های فارسی با توجه به جایگاه نویسه در واژه به صورت منفصل یا متصل نگاشته می‌شوند و این خود امر آموزش و همچنین ذخیره و بازیابی اطلاعات را دشوار می‌سازد. فرهنگستان زبان و ادب فارسی دستورالعملی را تحت عنوان دستور خط فارسی تدوین کرده است تا به یکسان‌سازی و بهبود نگارش فارسی منجر شود یکی از مباحث این دستورالعمل به فاصله‌گذاری مربوط است. طبق این دستورالعمل در خط فارسی از دو فاصله یعنی نیم فاصله (درون کلمه) و فاصله کامل (برون کلمه) استفاده می‌شود رعایت فاصله کامل و نیم فاصله در پرهیز از بدخوانی بسیار راه‌گشا است (فرهنگستان زبان و ادب فارسی 1389: 10). برای مثال در مورد نام و نام خانوادگی ایرانیان تنها فاصله‌گذاری دقیق است که باعث خوانایی و رفع ابهام می‌شود. زنجیره «علی رضا خانی» دو خوانش محتمل «علیرضا خانی» و «علی رضاخانی» دارد که تنها فاصله‌گذاری صحیح می‌تواند خوانش درست را تعیین نماید صفار مقدم (1386، 125) در بازیابی اطلاعات فاصله‌گذاری مخصوصاً در مورد واژه‌های مشتق و مرکب نقش تعیین‌کننده‌ای دارد چرا که پیوسته نویسی جدا نویسی کلید واژه‌های جستجو نتایج متفاوتی را بدست می‌دهد به عنوان نمونه پایگاه مگ ایران قابلیت جستجوی واژگان دو قسمتی را از طریق استفاده از نقطه فراهم آورده است اما کاربر با جستجوی واژه «خاکبرداری» در سه حالت «خاک برداری»، «خاکبرداری» و «خاک برداری» به نتایج متفاوتی دست خواهد یافت با جستجوی واژه «خاک برداری» 5 یافته بازیابی می‌شود در حالی که با جستجوی «خاک برداری» نتایج جستجو به 16 صفحه نیز می‌رسد. در میان این 16 صفحه یک نتیجه در برگیرنده واژه «خاک برداری» است در حالی که در میان نتایج 5 یافته اول نبوده و از 5 یافته اول حاصل از جستجوی «خاک برداری» مورد در نتایج حاصل از جستجوی «خاک برداری» مشاهده نشد. این در حالی است که نتایج حاصل از جستجوی «خاکبرداری» 7 مورد بوده که کاملاً متفاوت از دو جستجوی قبلی است. در واقع با جستجوی «خاک برداری» آن دسته از نتایج مرتبط در سایر جستجو‌ها به عنوان مرتبط ترین شناسایی نشده و در صفحات 2 و 4 و 8 پراکنده اند. در پایگاه RCeST نیز جستجوی کلیدواژه «خاک + برداری» 496، «خاکبرداری» 18 و عبارت «خاک برداری» 9 نتیجه را بازیابی می‌نماید بررسی نتایج جستجوهای «خاکبرداری» و «خاک برداری» و مقایسه آن‌ها با یکدیگر نشان داد که تنها یکی از عنوان‌ها مشترک بوده است. این بدان معنا است که در هر دو پایگاه

ص: 163

1- البته در نتایج بازیابی شده برای این عبارت، جستجو هم‌چنان مسائلی همچون عدم رعایت درست تقطیع به چشم می‌خورد

Letter -2

جستجو با نگرش های متفاوت به ریزش کاذب منتهی می شود. بدیهی است تعداد زیاد بازیافت ها و کاهش دقت در 10 نتیجه اول، یافتن نتایج مرتبط را برای کاربر زمان بر و دشوار خواهد کرد. باید توجه نمود که واژه هایی مانند خاک برداری کتاب، شناسی خودکشی و غیره از طریق فرایند ترکیب و اشتقاق ناپایگانی ساخته شده اند (شقایق 1386، 99) و اجزای دوم آن ها یعنی «-برداری»، «-شناسی» به طور مستقل بکار نرفته و استقلال واژگانی و معنایی ندارند بنابراین نظام های بازیابی اطلاعات باید راهکاری را برای جستجو در نظر گیرند تا این اجزا بصورت مجزا جستجو و بازیابی نشده و در ترکیبات مورد نظر به صورت یک واحد جستجو شوند.

مسائل صرفی و بازیابی اطلاعات

صرف (1) یکی از شاخه های زبان شناسی است که به مطالعه واژه و ساختمان درونی آن ها می پردازد. در این نوشتار به دو موضوع صرفی مرتبط به بازیابی اطلاعات یعنی واژه مرکب و وام گیری واژگانی پرداخته می شود. همان طور که پیش تر گفته شد واژه مرکب در فارسی چهار نوع مختلف برون مرکز، درون مرکز، دوسویه و متوازن دارد و ویژگی های خط فارسی در اغلب موارد باعث می شود که نتوان از لحاظ نوشتاری تمایزی بین واژه مرکب و گروه نحوی قائل شد. در نتیجه نظام های بازیابی اطلاعات نیز نمی توانند واژه مرکب را از گروه نحوی تشخیص دهند. برای مثال واژه «زیست شیمی» (2) یک واژه مرکب درون مرکز هسته پایانی است که «شیمی» هسته آن را تشکیل می دهد. در حالی که نظام بازیابی اطلاعات این واژه مرکب را هم گروه نحوی «زیست شیمی» و هم «زیست+ شیمی» قلمداد می کند. جستجوی این واژه مرکب در گوگل 11800000 بازیافت در پی دارد که بررسی 10 نتیجه اول نشان می دهد که نظام بازیابی هر سه حالت را مد نظر داشته است. این در حالی است که تنها 50000 بازیافت مرتبط محسوب می شود. لازم به ذکر است که در واژه مرکب «زیست شیمی»، «زیست» پیشوند است و نه واژه مستقل و «شیمی» هسته است در حالی که در واژه «محیط زیست» «زیست» واژه ای مستقل به شمار می رود و «محیط» هسته است. در واقع «زیست» در دو مثال فوق الذکر از دو مقوله و نقش متفاوت برخوردار است که اگر نظام های بازیابی مجهز به امکاناتی برای تشخیص و تعیین نوع آن ها بودند، مسلماً بازیابی بهبود چشم گیری می یافت. از دیگر چالش های صرفی بازیابی اطلاعات وام گیری واژگانی است. همان طور که می دانیم زبان ها از شیوه های مختلفی برای افزایش و غنی سازی واژگان خود بهره می برند که یکی از آن ها وام گیری است. در وام گیری، زبانی زبان مقصد عناصر زبانی را از زبان مبدأ به صورت های مختلفی مانند وام گیری، مستقیم غیر مستقیم ترجمه قرضی تعبیر، قرضی تغییر قرضی، ترجمه و تعبیر، قرضی آمیزش قرضی و تبادل قرضی وام می گیرد یکی از معضلات بازیابی اطلاعات خصوصاً در حوزه های تخصصی و علمی حاصل از وام گیری مستقیم است. در وام گیری مستقیم واژه ای از زبان

ص: 164

الف مستقیماً وارد زبان ب می شود این واژه پس از ورود به زبان ب بر طبق قواعد آوایی زبان مقصد تغییراتی پیدا می کند تا تلفظ آن برای گویشوران آسان شود (شقایقی 1386، 127-131). با توجه به این که در حوزه های تخصصی سرعت استفاده و جذب علم و یا فناوری بالاتر از سرعت واژه گزینی در زبان مقصد می باشد، لذا بسیاری از واژگان به همان صورت زبان مبدأ در زبان مقصد آوانگاری می شوند. بدیهی است چنین فرایندی عمیقاً تحت تأثیر نحوه تلفظ متخصص زمینه موضوعی بوده به طوری که امکان دارد یک واژه با املاهای متفاوتی نوشته شود؛ برای مثال سندروم / سندرم کلسیم / کالسیم، نیدروژن / هیدروژن این چندگانه نویسی باعث می شود که نه تنها ذخیره بلکه بازیابی اطلاعات نیز با چالش های فراوانی رو به رو شود به عنوان نمونه جستجوی معادل های آوانگاری شده واژه Psychology در پایگاه تخصصی مجلات نور ادعای پیش گفته را تأیید می کند برای واژه «پسیکولوژی» 95 یافته «پسیکولوژی» 3 یافته و «سایکولوژی» 8 یافته به دست آمد. لازم به ذکر است این نگارش های متفاوت به ترتیب بر اساس تلفظ، فرانسه آلمانی و انگلیسی وارد زبان فارسی شده است.

مسائل معنایی و بازیابی اطلاعات

زبان طبیعی از بُعد معنایی آکنده از ابهام است با وجود این کاربران زبان با استفاده از دانش زبانی و زمینه گفتمان در برخورد با این مسائل به ابهام زدایی می پردازند فلاحتی (1385) به طور کلی انواع ابهام های واژگانی را که می توانند به ابهام معنایی بی انجامند به 5 گروه عمده تقسیم می کند. ابهام مقوله ای ابهام حاصل از هم آوایی و هم نگاشتی ابهام چند معنایی و ابهام انتقالی ابهام مقوله ای در پیوند با معانی متفاوت یک واژه در بافت ها و نقش های مختلف همچون اسم، فعل، صفت و قید ایجاد می شود. در زبان فارسی این گونه از ابهام ها معمولاً به واسطه گذاشتن علائم واکه های کوتاه تا حد زیادی برطرف می شوند؛ اما مشکل اساسی این جاست که در متون مخصوص بزرگسالان و افراد باسواد (از جمله متون تخصصی و نوشتارهای علمی) استفاده از این علائم و مشخص کردن آن ها مرسوم نیست بنابراین در بسیاری از مواقع این خواننده است که با مراجعه به متن و بستر واژه، کار ابهام زدایی را انجام می دهد لازم به ذکر است در پاره ای از موارد امکان دارد تفاوتی میان تلفظ دو واژه وجود نداشته باشد اما همچنان به واسطه متفاوت بودن نقش دستوری هر واژه در جمله معنای متفاوتی از واژه مورد نظر دریافت گردد. به عنوان نمونه می توان به واژه «بردار» اشاره کرد.

بُردار: وزنه بردار دیابتی! قوی ترین مرد جهان (ورزش)

بُردار: پیشرفت هایی در کنترل بردار رانش (ریاضیات)

بُردار: داستان بردار کردن حسنگ وزیر (ادبیات)

چنان چه این واژه در بستر های مختلف جستجو شود با توجه به این که در متن اصلی استفاده از علائم نشان دهنده واکه های کوتاه مرسوم نیست و حتی در دو مورد اول و سوم نمایش واکه نیز تفاوتی در ابهام معنایی این تک واژه ایجاد نمی کند لذا در صورت جستجو، حداقل در سه گروه

متفاوت مقالات و مطالبی بازیابی می شود این در حالی است که فعل امر «بردار» از این گروه حذف شده است و حال آن که در جستجو و بازیابی اطلاعات احتمال بازیابی گروهی از مطالب که دارای این واژه هستند نیز وجود دارد تصاویر 1 (الف و ب) شمایی از نتایج جستجو در دو پایگاه اطلاعاتی فارسی را به نمایش می گذارد همان گونه که در این تصویر مشاهده می شود الگوریتم های بازیابی اطلاعات در نظام های جستجو و بازیابی اطلاعات در شرایط معمول با استفاده از میزان حضور واژه در متن و این که واژه در کجای متن باشد به رتبه بندی نتایج جستجو می پردازند. چنین روشی در پیوند با حالت های معمول زبانی مشکلی ایجاد نمی کند اما در مواجهه با مشکلات زبانی می تواند منجر به بروز ریزش کاذب در نتایج شود. لازم به توضیح است دو نمونه ارائه شده در تصویر الف هر دو برگرفته از صفحه اول نمایش نتایج هستند که انتظار می رود مرتبط ترین نتایج را در اختیار کاربر قرار دهند.

ابهام های حاصل از هم آوایی هم نگاشتی و چند معنایی گونه های دیگری از ابهام های زبان فارسی به شمار می آیند در این گروه از ابهام های واژگانی واژه ها یا تلفظ یکسانی دارند و یا به گونه یکسانی نوشته می شوند در نظام های بازیابی مبتنی بر جستجوی متنی ابهام های حاصل از هم نگاشتی مشکلاتی را برای جستجوگران اطلاعات به همراه دارد این مشکلات در پاره ای از اوقات با ابهام مقوله ای هم پوشانی پیدا می کنند یکی از نمونه های بارز این نوع از ابهام حداقل در پیوند با حوزه ای فیزیک واژه «گرم» و «گِرم» است (تصویر 2)



تصویر 1 (الف و ب). نتایج بازیابی شده از دو پایگاه مقالات فارسی مرکز منطقه ای اطلاع رسانی علوم و فناوری نشریات کشور (مگ ایران). در هر دو مورد بازیابی بر مبنای الگوریتم‌های معمول محاسبه ربط با توجه به حضور واژه در جمله انجام گرفته است.

تصویر 1 (الف و ب). نتایج بازیابی شده از دو پایگاه مقالات فارسی مرکز منطقه ای اطلاع رسانی علوم و فناوری و بانک اطلاعات نشریات کشور (مگ ایران). در هر دو مورد بازیابی بر مبنای الگوریتم‌های معمول محاسبه ربط با توجه به حضور واژه در جمله انجام گرفته است



تصویر ۲. نمایی از نتایج جستجو در پایگاه مجلات تخصصی نور

در همین راستا در حوزه معنایی^۱ چالش دیگری تحت عنوان هم‌معنایی مطرح می‌شود. هم‌معنایی به وضعیتی گفته می‌شود که دو یا چند واژه به یک مفهوم واحد اشاره داشته باشند (صفوی ۱۳۸۷، ۱۶۰). بایستی به این نکته توجه داشت که بسیاری از علوم در ایران بومی نبوده و پژوهشگران و دانشجویان غالباً از طریق ترجمه آثار غربی با آنها آشنا می‌شوند. بنابراین، اصطلاح‌گزینی برای واژگان علمی توسط یک مرجع واحد صورت نگرفته و گاه مترجمان و پژوهشگران با سلیقه شخصی دست به واژه‌گزینی می‌زنند. به همین دلیل با وجود اینکه یکی از مهم‌ترین عوامل موفقیت ارتباط علمی شفافیت واژگانی است، در بسیاری از زمینه‌های موضوعی برای بیان یک مفهوم واحد از چندین واژه متفاوت فارسی استفاده می‌شود. برای مثال در حوزه زبان‌شناسی برای هر مفهوم به طور متوسط سه واژه علمی فارسی وجود دارد (احمدی نسب ۱۳۹۰، ۲۷۶). به همین دلیل، همانگونه که پیشتر نیز بیان شد، بازایی اطلاعات در بسیاری از نظام‌ها برپایه واژگان موجود در متن صورت می‌پذیرد، و از آنجا که در متون مختلف مفاهیم به صورت‌های مختلف ثبت می‌شوند، لذا با جستجوی یک صورت از مفهوم، کاربر امکان دسترسی به صورت‌های دیگر همان مفهوم که منطقیاً مرتبط با نیاز اطلاعاتی وی خواهد بود را از دست می‌دهد. برای مثال برای مفهوم *Ontology* در فارسی سه معادل هستی‌شناسی، هستان‌شناسی و آنتولوژی بکار می‌رود که جستجوی هرکدام از این کلیدواژه‌ها تنها مدارکی را بازایی می‌کند که صرفاً همان صورت واژگانی در آنها بکار رفته است در حالی که هر سه واژه به یک مفهوم اشاره دارد (جدول ۱).

1. Synonymy

تصویر ۲. نمایی از نتایج جستجو در پایگاه مجلات تخصصی نور

در همین راستا در حوزه معنایی (1) چالش دیگری تحت عنوان هم‌معنایی مطرح می‌شود. هم‌معنایی به وضعیتی گفته می‌شود که دو یا چند واژه به یک مفهوم واحد اشاره داشته باشند صفوی (1387، 160) بایستی به این نکته توجه داشت که بسیاری از علوم در ایران بومی نبوده و پژوهش‌گران و دانشجویان غالباً از طریق ترجمه آثار غربی با آن‌ها آشنا می‌شوند بنابراین اصطلاح‌گزینی برای واژگان علمی توسط

یک مرجع واحد صورت نگرفته و گاه مترجمان و پژوهش‌گران با سلیقه شخصی دست به واژه‌گزینی می‌زنند به همین دلیل با وجود این که یکی از مهم‌ترین عوامل موفقیت ارتباط علمی شفافیت واژگانی است در بسیاری از زمینه‌های موضوعی برای بیان یک مفهوم واحد از چندین واژه متفاوت فارسی استفاده می‌شود. برای مثال در حوزه زبان‌شناسی برای هر مفهوم به طور متوسط سه واژه علمی فارسی وجود دارد (احمدی نسب 1390، 276). به همین دلیل همان‌گونه که پیش‌تر نیز بیان شد بازایی اطلاعات در بسیاری از نظام‌ها بر پایه واژگان موجود در متن صورت می‌پذیرد، و از آن‌جا که در متون مختلف مفاهیم به صورت‌های مختلف ثبت می‌شوند، لذا با جستجوی یک صورت از مفهوم کاربر امکان دسترسی به صورت‌های دیگر همان مفهوم که منطقاً مرتبط با نیاز اطلاعاتی وی خواهد بود را از دست می‌دهد. برای مثال برای مفهوم *Ontology* در فارسی سه معادل هستی‌شناسی، هستان‌شناسی و آنتولوژی بکار می‌رود که جستجوی هرکدام از این کلیدواژه‌ها تنها مدارکی را بازایی می‌کند که صرفاً همان صورت واژگانی در آن‌ها بکار رفته است در حالی که هر سه واژه به یک مفهوم اشاره دارد (جدول 1)

ص: 168

آسیب شناسی زبان و خط فارسی ... 169

جدول ۱. تفاوت بازیافت‌ها در پیوند با سه صورت واژگانی متفاوت از یک مفهوم واحد

بازگاه مقالات فارسی مرکز منطقه ای (Ricest)	بانک اطلاعات نشریات کنسور (مگ ایران)	بازگاه مجلات تخصصی نور
276	420	3877
8	7	69
8	5	100

ابهام انتقالی در بحث ترجمه از یک زبان به زبان دیگر مطرح می‌شود. زمانی که واژه از زبان مبدأ به زبان مقصد آمده اما در زبان مقصد چندین برابر نهاده پیدا می‌کند. برابرنهاده‌هایی که می‌توانند معانی متفاوت داشته و یا معانی مشابه و نگاشت‌های متفاوتی داشته باشند. این دست از واژگان ابهام‌برانگیز زمانی ایجاد مشکل می‌کنند که ابزار جستجو و بازیابی اطلاعات در قالب یک نظام بازیابی اطلاعات بین‌زبانی نیز فعالیت کند. به عنوان نمونه ۴ واژه *Transpiration* (تعرق)، *Sweating* (تعرق)، *Evapotranspiration* (تعرق و تبخیر)، و *Guttation* (تعریق) هر کدام در زبان انگلیسی به یک معنا به کار برده می‌شوند مثلاً در زبان مبدأ *Transpiration* برای گیاهان و در حوزه کشاورزی، و *Sweating* در پیوند با انسان به کار می‌رود. این درحالی است که واژه‌های پیش گفته در زبان فارسی معادل نسبتاً یکسانی دارند. هم اکنون گوگل تاحدی قابلیت جستجوی بین‌زبانی در سطح وب را فراهم آورده است. اما زمانی که معادل واژه *Transpiration* در صفحات فارسی جستجو شود نتایج مشابه با تصویر ۳ بدست خواهد آمد. این درحالی است که نتایج زمانی که واژه *تعرق* به فارسی برای بازیابی صفحات انگلیسی به کار رود به مراتب از پراکندگی بیشتری برخوردار خواهد بود؛ چرا که واژه *تعرق* در زبان فارسی برای بیان سه مفهوم متفاوت در زبان انگلیسی انتخاب شده است.



تصویر ۳. نمایش از صفحه اول نتایج بازیابی شده بر اساس اولین برابرنهاده پیشنهادی گوگل برای واژه

Transpiration

جدول 1 تفاوت بازیافت‌ها در پیوند با سه صورت واژگانی متفاوت از یک مفهوم واحد

ابهام انتقالی در بحث ترجمه از یک زبان به زبان دیگر مطرح می‌شود. زمانی که واژه از زبان مبدأ به زبان مقصد آمده اما در زبان مقصد چندین برابر نهاده پیدا می‌کند برابر نهاده‌هایی که می‌توانند معانی متفاوت داشته و یا معانی مشابه و نگاشت‌های متفاوتی داشته باشند. این دست از واژگان ابهام‌برانگیز زمانی ایجاد مشکل می‌کنند که ابزار جستجو و بازیابی اطلاعات در قالب یک نظام بازیابی اطلاعات بین

زبانی نیز فعالیت . کند به عنوان نمونه 4 واژه Transpiration (تعرق)، Sweating (تعرق)، Evapotranspiration (تعرق و تبخیر) و Guttation (تعریق) هر کدام در زبان انگلیسی به یک معنا به کار برده می شوند مثلاً در زبان مبدأ Transpiration برای گیاهان و در حوزه کشاورزی، و Sweating در پیوند با انسان به کار می رود این در حالی است که واژه های پیش گفته در زبان فارسی معادل نسبتاً یکسانی دارند هم اکنون گوگل تا حدی قابلیت جستجوی بین زبانی در سطح وب را فراهم آورده است. اما زمانی که معادل واژه Transpiration در صفحات فارسی جستجو شود نتایجی مشابه با تصویر 3 بدست خواهد آمد. این در حالی است که نتایج زمانی که واژه تعرق به فارسی برای بازیابی صفحات انگلیسی به کار رود به مراتب از پراکندگی بیش تری برخوردار خواهد بود؛ چرا که واژه تعرق در زبان فارسی برای بیان سه مفهوم متفاوت در زبان انگلیسی انتخاب شده است.

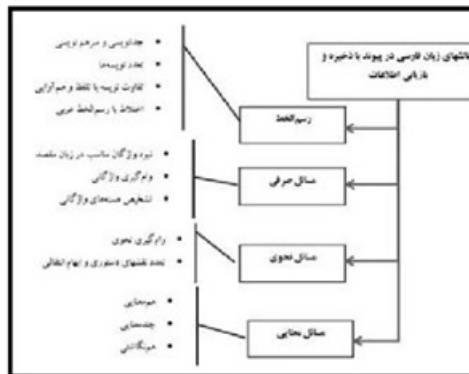
تصویر 3 نمایی از صفحه اول نتایج بازیابی شده بر اساس اولین برابر نهاده پیشنهادی گوگل برای واژه Transpiration

ص: 169

مسئله ایجاد یک نظام برای به هم نزدیک سازی ذهنیت پدیدآورنده و استفاده کننده از مسائل مهم در امر طراحی پایگاه های اطلاعاتی است. دسترسی مناسب به معنای ایجاد شرایط لازم از نظر زبان و واژگان اختصاص یافته به مدارک است؛ واژگانی که به کاربر در کاوش منابع یاری می رساند (مک ایوان 1379، 4) حال مشکل از این جا آغاز می شود که در پاره ای از پایگاه های اطلاعاتی، شخص سومی تحت عنوان نمایه ساز وارد شده و پاره ای از کلیدواژه ها را به مدارک با هدف بازیابی بهتر اختصاص می دهد در این فرایند با توجه به این که حلقه سومی در زنجیره ارتباطی میان پدیدآور و مخاطب ایجاد می شود اگر چه هدف بهبود بازیابی است اما گاه احتمال افزایش ضریب خطا اگر چه در ورود اطلاعات و چه در اختصاص واژگان کلیدی نیز وجود دارد (شاپوری 1379؛ Gross Taylor 2005) مسلماً دانش حوزه ای از طریق مطالعه حوزه ای حاصل می شود تحقیقات متعددی در تأیید این مطلب صورت گرفته است که وایلدرموت (Wildermuth، 2004) ضمن برشمردن برخی از آن ها تأکید می کند دانش حوزه ای و واژگانی افراد که عموماً برگرفته از متون تخصصی، است نقش عمده ای در انتخاب واژگان جستجو دارد در چنین شرایطی وظیفه نمایه ساز به عنوان حلقه واسط نزدیک کردن دو زبان پدیدآور و مخاطب است خصوصاً در هنگامی که نظام بازیابی اطلاعات قابلیت برقراری این ارتباط را بصورت مؤثر ندارد با وجود، این حجم زیاد متون ذخیره شده و کمبود نیروی انسانی توانمند متخصصان حوزه بازیابی اطلاعات را به سمت هرچه هوشمند ساختن بازیابی اطلاعات سوق داده است. تصویر 4 به منظور جمع بندی ساختواره ای درختی از اهم چالش های متصور در زمینه ذخیره و بازیابی اطلاعات را به نمایش می گذارد

عکس

مسئله ایجاد یک نظام برای به هم نزدیک‌سازی ذهنیت پدیدآورنده و استفاده‌کننده از مسائل مهم در امر طراحی پایگاه‌های اطلاعاتی است. دسترسی مناسب، به معنای ایجاد شرایط لازم از نظر زبان و واژگان اختصاص یافته به مدارک است؛ واژگانی که به کاربر در کاوش منابع یاری می‌رساند (سکایوان ۱۳۷۹، ۴۹). حال، مشکل از این جا آغاز می‌شود که در پاره‌ای از پایگاه‌های اطلاعاتی، شخص سومی تحت عنوان نمایه‌ساز وارد شده و پاره‌ای از کلیدواژه‌ها را به مدارک با هدف بازیابی بهتر اختصاص می‌دهد. در این فرایند با توجه به اینکه حلقه سومی در زنجیره ارتباطی میان پدیدآور و مخاطب ایجاد می‌شود، اگرچه هدف بهبود بازیابی است، اما گاه احتمال افزایش ضریب خطا چه در ورود اطلاعات و چه در اختصاص واژگان کلیدی نیز وجود دارد (شاپوری ۱۳۷۹؛ Gross & Taylor ۲۰۰۵). مسلماً دانش حوزه‌ای از طریق مطالعه حوزه‌ای حاصل می‌شود. تحقیقات متعددی در تأیید این مطلب صورت گرفته است که وایلدروموت (Wildermuth, ۲۰۰۴) ضمن برشمردن برخی از آنها، تأکید می‌کند دانش حوزه‌ای و واژگانی افراد که عموماً برگرفته از متون تخصصی است، نقش عمده‌ای در انتخاب واژگان جستجو دارد. در چنین شرایطی وظیفه نمایه‌ساز به عنوان حلقه واسط، نزدیک کردن دو زبان پدیدآور و مخاطب است. خصوصاً در هنگامی که نظام بازیابی اطلاعات قابلیت برقراری این ارتباط را بصورت مؤثر ندارد. با وجود این، حجم زیاد متون ذخیره شده و کمبود نیروی انسانی توانمند، متخصصان حوزه بازیابی اطلاعات را به سمت هرچه هوشمند ساختن بازیابی اطلاعات سوق داده است. تصویر ۴ به منظور جمع‌بندی ساختارهای درختی از اهم چالش‌های متصور در زمینه ذخیره و بازیابی اطلاعات را به نمایش می‌گذارد.



تصویر ۴. ساختار درختی چالش‌های عمده زبانی در پیوند با ذخیره و بازیابی اطلاعات^۱

۱. لازم به توضیح است این ساختار کامل‌تر از بحث‌های مطرح در این نوشتار است. دلیل این امر لزوم توجه به برخی از مهمترین موارد به منظور جلوگیری از طولانی شدن مبحث بوده است.

تصویر 4. ساختار درختی چالش‌های عمده زبانی در پیوند با ذخیره و بازیابی اطلاعات (1)

ص: 170

1- لازم به توضیح است این ساختار کامل‌تر از بحث‌های مطرح در این نوشتار است. دلیل این امر لزوم توجه به برخی از مهم‌ترین موارد به منظور جلوگیری از طولانی شدن مبحث بوده است

راهبردهای متفاوتی در بحث رفع این مسائل مطرح شده است که می توان آن ها را به گروه های مختلف دسته بندی نمود از میان آن ها عمده ترین موارد شامل تدوین پیکره زبان فارسی یادگیری ماشینی ریشه یاب ها، نرمال سازی، نویسه ها استفاده از الگوریتم های احتمالاتی همچون N-gram ها و مهار واژگانی می شود برای رفع چالش های معنایی بازیابی اطلاعات به نظر می رسد که باید بر اساس یک پیکره زبانی و استخراج قواعد واژه سازی زبان فارسی و تعیین اجزای واژگانی همچون وندها و ریشه ها عمل کرد این کار از طریق یادگیری ماشینی از یکسو و ایجاد پیکره های زبان فارسی به صورت الکترونیک امکان پذیر است حرکت هایی در این زمینه در پیوند با زبان فارسی صورت گرفته است اما همچنان در مراحل ابتدایی است در این راستا یکی از کارهای ارزشمند انجام شده پارس مورف متعلق به مواجی، اسلامی و وزیر نژاد (1390) است. این نرم افزار قابلیت تحلیل صرفی زبان فارسی را داراست. در همین رابطه می توان به ریشه یاب تهیه شده توسط احسان و فیلی (1390) نیز اشاره کرد که میزان دقت در نظام های بازیابی را 5 درصد افزایش می دهد. از دیگر تلاش ها می توان به ریشه یاب برنجیان (1390) برای افعال ماضی و مضارع ناگذر اشاره کرد. لازم به یادآوری است که حرکت های مشابه برای زبان عربی مدت هاست آغاز شده و به مرحله ارزیابی و بهسازی رسیده است. به عنوان نمونه می توان به پژوهش های محمود ماجد و خلدون (Mahmoud, Majed Kaldoun 2011) و یا انیس و دیگران (Anis et al. In Press) اشاره کرد.

از راهبرد های پیشنهادی در پیوند با برطرف سازی مشکلات مرتبط با تعدد نویسه هایی همچون الف و یا همزه در زبان عربی نرمال سازی است. در این راهبرد با استفاده از یک الگوریتم ساده، کلیه صورت های مختلف یک نویسه به طور خودکار به یک شکل یکسان تبدیل می شود به عنوان مثال کلیه نگارش های متفاوت «آ» و «ا» به صورت «1» در نظر گرفته می شود (Mahmoud, Majed Kaldoun 2011). اما با بررسی های صورت گرفته و همان گونه که پیش تر در این مقاله نیز مطرح گردید زبان فارس به واسطه ویژگی های معنایی خاص خود بر تابنده چنین تغییری، حداقل در پیوند با این مثال نمی باشد لذا نرمال سازی در پیوند با نویسه های زبان فارسی بایستی با دقت و مطالعه بیش تری صورت گرفته و نمی توان به یافته های مطالعات مرتبط با زبان عربی در این زمینه استناد کرد راهبرد دیگر در این راستا استفاده از ریشه یاب ها است استفاده از ریشه یاب امکان بازیابی صورت های مختلف واژگانی که به یک حوزه مفهومی واحد اشاره دارند را فراهم می کند. این در حالی است که این راهبرد می تواند به صورت بالقوه پاسخی به مشکلات نرم افزارهای بازیابی اطلاعات در پیوند با صورت های مفرد و جمع واژگان خصوصاً جمع های مکسر و یا واژه های هم خانواده باشد (Mahmoud, Majed Kaldoun 2011) با توجه به این که زبان فارسی از نظر رسم الخط مشابه با زبان عربی است پیشنهاد می شود در این زمینه از مطالعات صورت گرفته در پیوند با خط عربی و بازیابی اطلاعات استفاده شود در حالی که در زمینه ویژگی های صرفی و معنایی زبان فارسی به زبان انگلیسی نیز نزدیک بوده و مطالعات صورت گرفته در این زمینه تا حدود زیادی

اگر چه در بسیاری از متون بحث رعایت فاصله گذاری مناسب به پدیدآور نسبت داده می شود اما هم چنان نیازمند نظام هایی با قابلیت تشخیص مناسب فاصله ها در واژگان مرکب (1) هستیم. در زبان فارسی یکی از شیوه های پر کاربرد واژه سازی استفاده از فرایند ترکیب است لذا تعداد واژگان ترکیبی در زبان فارسی بسیار زیاد است این در حالی است که در زبان فارسی به واسطه ماهیت زبانی و دستوری و نیز سهولت، خواندن واژگان مرکب معمولاً به صورت ناپیوسته همراه با یک فاصله یک حرفی ثبت می شوند. لازم به یادآوری نیست که در زبان انگلیسی این مشکل کمتر دیده می شود. به عنوان نمونه کلیه واژگانی که با «شناسی» همراه هستند. معمولاً به صورت دو جزء جداگانه ثبت می شوند در حالی که در زبان انگلیسی "logy" بخشی از واژه به شمار می آید. بدیهی است رفع این مشکل تنها از طریق استفاده از عملگر And در یک فیلد جستجو امکان پذیر نمی باشد، چرا که منجر به افزایش تعداد بازیافت های غیر مرتبط می شود این در حالی است که چنان چه کاربری کلمات وارد شده در یک فیلد به شکل عبارتی تغییر یابد امکان بروز این خطا کمتر خواهد شد. اما با نگاهی به متون موجود در پیوند با رفع این دسته از چالش ها، حداقل در پیوند با زبان عربی استفاده از الگوریتم های N-gram پیشنهاد می شود. دلیل این امر قابلیت ارائه بافت واژه مورد جستجو همراه با واژه در قالب پیشنهادهایی برای بهسازی جستجوست (Mahmoud, Majed Kaldoun 2011).

مهار واژگانی یکی دیگر از راهبرد های مفید در زمینه بهسازی جستجو از طریق کم اثر سازی تعدد صورت های واژگانی و یا معنایی است بدین معنا که با ارائه صورت های مرجح، نامرجح و واژگان شامل و زیر شمول کلیه صورت های متفاوت یک مفهوم را ذیل یک مفهوم واژه پذیرفته شده گرد می آورد. واژگان مهار شده معمولاً نمایان گر ساختواره ای از روابط معنایی و سلسله مراتبی بوده و استفاده از آن خصوصاً در برطرف ساختن چالش های هم معنایی، چندمعنایی، املاهای متعدد و وام گیری واژگانی کاراست (Svenonius 2003). این در حالی است که استفاده همزمان از اصطلاحنامه ها در کنار کلید واژه های زبان طبیعی یکی از کاراترین روش ها در این راستا به شمار می آید که از جمله پایگاه های موفق در این زمینه می توان به INSPEC، Eric و LISA اشاره کرد.

در پایان خاطر نشان می سازد که تاکنون پیکره های فارسی مختلفی همچون همشهری محک بی جن خان و دادگان زبان فارسی تدوین شده است از این میان درودی و دیگران (1387) مجموعه محک وب را تهیه کرده اند که از آن می توان برای انجام مطالعات در حوزه بازیابی اطلاعات فارسی در وب استفاده نمود. این مجموعه dotIR نام دارد و استفاده از آن برای امور غیر تجاری رایگان است. مجموعه پیش گفته از یک پیکره استاندارد 50 پرس و جوی استاندارد 18 هزار دآوری تعیین ربط پرس و جوها به اسناد پیکره و 50 هزار بردار ویژگی استخراج شده از اسناد، تشکیل شده و بیش تر مناسب پژوهش در رابطه با بازیابی اطلاعات است همشهری (2) مجموعه دیگری است که

1- در این جا برای سهولت بحث واژه مرکب به واژه غیر بسیط اطلاق شده و معنای دقیق زبان شناختی آن مدنظر نیست 2.

<http://ece.ut.ac.ir/dbrg/hamshahri/faindex.html>

2- <http://ece.ut.ac.ir/dbrg/hamshahri/faindex.html>

توسط رهگذر و دیگران بر اساس اخبار روزنامه همشهری تهیه شده و دارای دو نسخه 1 و 2 است که نسخه دوم از امکانات بیش تری مانند پیوند به تصاویر و به اصل صفحات وب برخوردار است. از این مجموعه نیز می توان برای انجام پژوهش در حوزه بازیابی اطلاعات، فارسی، پردازش زبان طبیعی و تدوین الگوریتم های ریشه یابی استفاده نمود. بر اساس این مجموعه پژوهش های متعددی صورت گرفته که از جمله می توان به آل احمد و دیگران (2007) اشاره نمود که با استفاده از این مجموعه مدل فضای برداری بر اساس n-gram و کلمه را ارزیابی کرده و نشان داده اند که بازیابی متون فارسی بر اساس مدل فضای برداری 4-gram و طرح وزن دهی Inu ltu نتایج قابل قبول و تحلیل محلی متن (1) بهترین نتایج را در بر خواهد داشت. از پیکره های زبانی دیگر می توان به پیکره بی جن خان (2) اشاره کرد که از متون خبری و متون عمومی شامل 4300 موضوع تشکیل شده و بر اساس 40 مقوله دستوری فارسی بر چسب دهی شده است و مناسب پژوهش های مرتبط به پردازش زبان طبیعی است. نمونه دیگر پایگاه دادگان فارسی (3) است که توسط عاصی و دیگران در پژوهش گاه علوم انسانی و مطالعات فرهنگی تهیه شده است این پایگاه گونه های مختلف زبانی نوشتاری و گفتاری را در بر می گیرد و در آن امکان جستجوی واژه ها ترکیب ها و بررسی بسامد آن ها به همراه گزارش های آماری از متون وجود دارد. در این جا ذکر این نکته خالی از لطف نیست که متون موجود و پژوهش های انجام شده نشان دهنده این امر است که علی رغم وجود این پیکره های زبانی فارسی، پژوهش های چندانی بر اساس این پیکره ها صورت نگرفته است شاید بتوان این امر را به عدم اطلاع اکثر پژوهش گران از اهمیت و نقش پیکره های زبانی در تحقیقات در حوزه پردازش زبان، طبیعی بازیابی اطلاعات و تهیه ریشه یاب ها و همچنین عدم آموزش کافی جامعه علمی در این زمینه نسبت داد. در این نوشتار نیز با توجه به چارچوب زمانی در نظر گرفته شده برای بررسی، حاضر متأسفانه محققان قادر به بررسی دقیق این پیکره ها و استفاده از آن ها در این بررسی نشدند. بنابراین لازم است تا پژوهشی جامع، میزان موفقیت راهبردهای پیشنهادی این بررسی بر پایه پیکره های پیش گفته را بسنجد. هم چنین توصیه می شود که در طراحی و به سازی پایگاه های اطلاعاتی زبان فارسی کارگروه هایی متشکل از متخصصان علم اطلاعات و دانش شناسی زبان شناسی و علوم رایانه نقش اصلی را ایفا نمایند.

منابع

آشوری، داریوش 1375 چند پیشنهاد دیگر برای اصلاح خط فارسی سرهم نویسی و جدا نویسی. نگاه نو. 28:101-117.

احمدی نسب، فاطمه. 1390 تهیه و تدوین اصطلاحنامه تک زبانه فارسی زبان شناسی رساله دکترا دانشگاه علامه طباطبایی.

اسلامی محرم 1381 دشواری های پردازش رایانه ای خط فارسی نشر دانش دوره 19(3) 28-32

ص: 173

Local context analysis -1

<http://ece.ut.ac.ir/DBRG/Bijankhan> -2

<http://ece.ut.ac.ir/DBRG/Bijankhan> -3

اسلامی، محرم 1386 خط فارسی و رسانه های گروهی دو فصل نامه پردازش علائم و داده ها. 8(2):93-98.

برنجیان شاپور رضا 1390. ریشه یاب ماضی و مضارع از مصدر افعال ناگذر در زبان فارسی. شیراز نوید شیراز مرکز منطقه ای اطلاع رسانی علوم و فناوری

دبیر مقدم، محمد. 1383. زبان شناسی نظری پیدایش و تکوین دستور زایشی (ویراست دوم) تهران انتشارات سمت.

درودی احسان و دیگران 1387 پیکره محک استاندارد برای تحقیقات بازایی اطلاعات وب فارسی گزارش، فنی گروه تحقیقاتی پایگاه داده ها دانشگاه تهران شماره: <http://ccc.ut.ac.ir/DBRG/webir/files/Papers/WebIR.pdf>. DBRG-TR-138702 (دسترسی در 1391/10/19).

شاپوری، سودابه (1379) بررسی مشکلات جستجوی موضوعی استفاده کنندگان رایانه ای کتابخانه مرکزی دانشگاه فردوسی مشهد پایان نامه کارشناسی ارشد مشهد دانشکده علوم تربیتی و روانشناسی

شقایقی ویدا. 1386. مبانی صرف تهران: انتشارات سمت.

صفر مقدم، احمد. 1386. فاصله گذاری در خط فارسی نامه فرهنگستان 9/4: 123-137.

صفوی، کورش 1387. درآمدی بر معنی شناسی. تهران: انتشارات سوره مهر (حوزه هنری سازمان تبلیغات اسلامی)

فرهنگستان زبان و ادب فارسی 1389. دستور خط فارسی چاپ نهم تهران: فرهنگستان ادب و زبان فارسی (نشر آثار).

فلاحی فومنی، محمدرضا 1385 ابهام در ماشین ترجمه کتابداری و اطلاع رسانی. 9(3): 21 - 38

کامری، برنارد جان ماونتفورد ویوین، لا و د. آ. کروز 1990. زبان های دنیا چهار مقاله در زبان شناسی ترجمه کورش صفوی 1384 تهران: انتشارات سعادت.

گل تاجی، مرضیه و سعیده برزگر. 1389 بررسی مشکلات ریخت شناسی زبان فارسی در سه پایگاه اطلاعاتی مرکز منطقه ای اطلاع رسانی علوم و فناوری پژوهشگاه اطلاعات و مدارک علمی ایران و جهاد دانشگاهی کتابداری و اطلاع رسانی 13(2): 191-214

مک ایوان اندرو. (1379) استفاده از سر عنوان های موضوعی کتابخانه کنگره: هزینه همکاری برای رسیدن به دسترس پذیری ترجمه مجتبی اسدی گزیده مقالات ایفلا 98 (آمستردام: 21 - 16 اوت 1998) (ص. 59 - 47) تهران کتابخانه ملی جمهوری اسلامی ایران.

محقق زاده محمد صادق و کاظم زارعیان 1383 ارائه راه حل برای برخی مسائل اتوماسیون و نگارش فارسی. فصلنامه اطلاع-رسانی 19(3 و 4): 1-10.

معصومی همدانی حسین 1381 خط فارسی و رایانه نشر دانش. 19(2) 2-6

مواجی وحید محرم اسلامی و بهرام وزیر نژاد 1390. پارس مورف: تحلیل گرافیکی زبان

- Aleahmad, Abolfazl. et al. 2007. N-Gram and Local Context Analysis for Persian Text Retrieval, International Symposium on Signal Processing and Its Applications, Sharjah U.A.E. Retrieved 2013-01-8
From <http://ece.ut.ac.ir/dbrg/hamshahri/files/Papers/isspa.pdf>
- Anis, Z. et al. In Press. Contribution to Semantic Analysis of Arabic language. Advances in Artificial Intelligence. Retrieved 2012-11-10 From <http://www.hindawi.com/journals/aai/aip/620461>
- .Crystal, David. 2008. A dictionary of linguistics and phonetics. Six Edition. Blackwell Publishing
- Gross, Tina, Arlene G. Taylor. 2000. What Have We Got to Lose? The Effect of Controlled Vocabulary on Keyword Searching Results. College Research Libraries. 66(3): 212 -230. Retrieved August 2, 2006, From www.ala.org/ala/acrl/acrlpubs/crljournal/backissues2005a/crlmay05/Gross.pdf
- Mahmoud, R., S. Majed, and Z. Kaldoun. 2011. Improving Arabic information retrieval system using N-gram method. WSEAS Transactions on Computers 10 (4): 125-133
- Svenonius, E. 2003. Design of Controlled Vocabularies. Encyclopedia of Library and Information Science) (822-838). New York: Marcel Dekker (reprinted from the first edition, 1989
- Wildermuth, Barbara M. 2004. The Effects of Domain Knowledge on Search Tactic Formulation. Journal of the American Society for Information Science and Technology. 55(3): 246-258

هدف: هدف اصلی از انجام این پژوهش ارزیابی کاربرد پذیری وب سایت نهاد کتابخانه های عمومی کشور است.

روش شناسی: پژوهش حاضر که از نوع کاربردی است، با استفاده از شیوه ارزیابانه به بررسی و ارزیابی معیارها و مولفه های مطرح در کاربرد پذیری وب سایت نهاد کتابخانه های عمومی کشور می پردازد. ابزار گردآوری اطلاعات در این پژوهش سیاهه ارزیابی محقق ساخته مشتمل بر 11 معیار و 160 مؤلفه است که وب سایت مورد نظر با آن سنجیده شد. در سیاهه مزبور از روش روایی صوری و محتوایی استفاده شده منظور تجزیه و تحلیل یافته های پژوهش از آمار توصیفی فراوانی درصد و میانگین برای به توصیف کشیدن وضعیت موجود وب سایت مورد بررسی استفاده شد.

یافته ها: در مجموع نتایج به دست آمده نشان داد در وب سایت نهاد کتابخانه های عمومی کشور نیمی از استانداردها 346 امتیاز (52/2 درصد) از 663 امتیاز وب سایت شاخص رعایت گردیده است.

کلید واژه ها: ارزیابی، کاربرد پذیری وب سایت، نهاد کتابخانه های عمومی کشور

صدیقه محمد اسماعیل (1) | ماهرخ ناصحی اسکویی (2)

مقدمه

امروزه از وب به عنوان یکی از مهم ترین و اصلی ترین ابزارها برای دسترسی به اطلاعات استفاده می کنند لذا ضروری است، کاربران وب میزان کاربرد پذیری منابع موجود بر روی آن (به ویژه وب سایت های آن) را مورد بررسی و ارزیابی مستمر قرار دهند در این راستا پژوهش حاضر بر آن است تا وب سایت نهاد کتابخانه های عمومی کشور را مورد ارزیابی قرار دهد تا از این طریق مشخص گردد وب سایت نهاد به لحاظ کاربردپذیر بودن در چه شرایطی قرار دارد و تا چه حد در طراحی آن اصول و ضوابط طراحی وب سایت به لحاظ رعایت مؤلفه های مربوط به این امر (اعم از: اعتبار، صحت، روزآمد بودن، سطح پوشش و مخاطبان خاص، وجود نماهای تعاملی و تبادل عینیت، اطلاعات قابلیت ناوبری، نماهای غیر متنی، دسترس پذیری، کارآمدی، ویژگی های ظاهری رعایت گردیده است. بدیهی است، پرداختن به وب سایت نهاد کتابخانه های عمومی کشور با عنایت به جایگاه محوری آن در بحث اطلاعات و اطلاع رسانی امری حائز ارزش و توجه است.

ص: 177

1- استادیار گروه کتابداری و اطلاع رسانی دانشگاه آزاد اسلامی واحد علوم تحقیقات تهران M.esmaeili2@gmail.com

2- دانشجوی کارشناسی ارشد کتابداری و اطلاع رسانی دانشگاه آزاد اسلامی واحد علوم و تحقیقات تهران

mahrokh.nasehi@gmail.com

2- روش پژوهش و توجیه روایی آن

روش پژوهش حاضر، روش کتابخانه‌ای (سندی) و پیمایشی از نوع ارزیابانه است. بدین معنا که، با استفاده از سیاهه کنترل یا به ارزیابی کاربرپذیری وب سایت نهاد کتابخانه‌های عمومی کشور (بر مبنای معیارهای متعدد موجود در چک لیست) پرداخته شده است. در توجیه روایی استفاده از چنین روشی جهت انجام این پژوهش ذکر این نکته ضروریست که از آن جایی که شاخص‌های بیرونی متعددی در قالب سیاهه واریسی بوده که چگونگی کاربرپذیری صفحات وب نهاد کتابخانه‌های عمومی کشور با آن‌ها سنجیده شده است لذا مناسب‌ترین راه برای آگاهی از وضعیت این وب سایت‌ها (بر اساس معیارها و مؤلفه‌های ذکر شده) روش ارزیابانه فوق بوده است.

3- شیوه گردآوری اطلاعات و تجزیه و تحلیل آن‌ها

در این پژوهش گردآوری اطلاعات با استفاده از روش مشاهده مستقیم وب سایت و بر اساس سیاهه‌ای با 11 معیار و 160 مؤلفه صورت گرفت و بر این مبنای وب سایت نهاد کتابخانه‌های عمومی کشور از نقطه نظر کاربردپذیری و معیارهای مطرح در آن ارزیابی گردید. برای این منظور، تلاش شد تا ابتدا نسبت به تهیه سیاهه مربوطه تحت عنوان سیاهه ارزیابی کاربرپذیری وب سایت (1) اقدام گردد. از این رو، متون و منابع و پیشینه‌های موجود در زمینه کاربرپذیری وب سایت‌ها در داخل و خارج کشور مورد مطالعه و تفحص قرار گرفت و به منظور اطمینان از روایی محتوایی هر چه بیش‌تر در اختیار بررسی 5 تن از مدیران وب سایت‌های کتابخانه‌ای قرار گرفت سپس نظرات آن‌ها، بررسی و نسبت به تهیه نسخه نهایی سیاهه (مشتمل بر 11 معیار کلی و 160 مؤلفه) اقدام شد و با استفاده از آن وضعیت وب سایت نهاد کتابخانه‌های کشور در طول مدت زمان انجام این پژوهش سه ماهه اول (1391)، مورد بررسی و ارزیابی قرار گرفت. به منظور تجزیه و تحلیل یافته‌های پژوهش از آمار توصیفی (فراوانی درصد و میانگین) برای به توصیف کشیدن وضعیت موجود استفاده شد. در سیاهه یاد شده از دو مقیاس وجود و نبود بلی (i) و خیر (-) استفاده شده که امتیازات در نظر گرفته شده برای آن‌ها به ترتیب عبارتند از بلی = 1 (یک) و خیر = 0 (صفر). (نکته: اعداد مندرج در قسمت تواتر در جدول به شماره منابع ارجاع خورده است).

4- تجزیه و تحلیل داده‌ها و ارائه یافته‌ها

همان‌گونه که قبلاً نیز گفته شد در این پژوهش از یک وب سایت مفروض استفاده شده با این فرض که در طراحی آن هر 160 مؤلفه یا ویژگی به گونه‌ای رعایت شده است که می‌توان آن را به مثابه شاخصی برای ارزیابی و سنجش دیگر وب سایت جامعه مورد مطالعه به کار برد. یافته‌های پژوهش بیان می‌کند که در وب سایت شاخص امتیاز کل رعایت همه معیارهای مطرح در کاربردپذیری سایت برابر با 663 امتیاز است که از این امتیاز 107 امتیاز مربوط به رعایت مؤلفه‌های اعتبار اطلاعات (16 درصد)

ص: 178

32 امتیاز مربوط به صحت اطلاعات (5 درصد)، 26 امتیاز مربوط به روزآمد بودن (4 درصد)، 21 امتیاز مربوط به سطح پوشش و مخاطبان خاص (3 درصد) 31 امتیاز مربوط به نماهای تعاملی و تبادلی (5 درصد)، 34 امتیاز مربوط به عینیت اطلاعات (5 درصد)، 166 امتیاز مربوط به رعایت معیار کلی ناوبری در سایت (جمعا: 25 درصد)، به تفکیک شامل: 9 امتیاز مرتبط با ویژگیهای عنوان مرورگر (1 درصد) 22 امتیاز مرتبط با ویژگیهای عنوان صفحه (3 درصد) 92 امتیاز مرتبط با پیوندهای متنی و فرامتنی (14 درصد)، 6 امتیاز مرتبط با نشانه اینترنتی (1 درصد)، 25 امتیاز مرتبط با نقشه سایت یا نمایه (4 درصد) 12 امتیاز مرتبط با موتور جستجوی داخلی (2 درصد)، 64 امتیاز مربوط به نماهای غیرمتنی (10 درصد)، 77 امتیاز مربوط به قابلیت دسترس پذیری (12 درصد)، 99 امتیاز مربوط به کارآمد پذیری (15 درصد) و 6 امتیاز مربوط به ویژگی های ظاهری (1 درصد) است

همان گونه که از جدول 1 مستفاد می شود یافته های تحقیق نشان می دهد که وب سایت نهاد کتابخانه های عمومی کشور از لحاظ رعایت معیارهای «اعتبار اطلاعات» 45 امتیاز (42 درصد)، «صحت اطلاعات» امتیاز (53 درصد)، «روزآمد بودن» 7 امتیاز (26 درصد)، «سطح پوشش و مخاطبین خاص» 10 امتیاز (47 درصد) «نماهای تعاملی و تبادلی» دارای 12 امتیاز (38 درصد)، «عینیت اطلاعات» 15 امتیاز (44 درصد) از لحاظ معیارهای کلی «ناوبری» 103 امتیاز (62 درصد) «نماهای غیر متنی» 18 امتیاز (28 درصد) «دسترس پذیری» 48 امتیاز (62 درصد)، «کارآمدی» 51 امتیاز (51 درصد) و در نهایت در رابطه با آخرین معیار مورد بررسی که ویژگی های «ظاهری» 6 امتیاز (1000 درصد) را به طور کامل به خود اختصاص داده است. یافته های این پژوهش با عصاره که به بررسی وب سایت های 58 کتابخانه ملی جهان با هدف شناسایی ویژگی های شاخص در طراحی وب سایت یک کتابخانه ملی به روش تحلیل محتوا و ارزیابی بر مبنای سیاهه واری (وب سایت کتابخانه ملی آمریکا 80 درصد، وب سایت کتابخانه بیمارستانی انگلستان 28 درصد و وب سایت کتابخانه ملی کشورهای مالزی و دانمارک 72 درصد با سیاهه واری و وب سایت کتابخانه ملی ایران 51 درصد با سیاهه همخوانی داشتند) و نیز با یافته های تحقیق حاجی زین العابدینی که به بررسی وضعیت وب سایت های کتابخانه های ملی جهان از نظر میزان مورد استفاده قرار گرفتن از طریق پیوندهای آن ها می پردازد همسو است. و در بعد دسترس پذیری نیز با یافته های محمد اسماعیل و کاظمی در تحقیقی با عنوان دسترس پذیری وب سایت کتابخانه های ملی، خاورمیانه با استفاده از سیاهه واری ای مشتمل بر 13 مؤلفه که در آن به بررسی تطبیقی وب سایت سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران با وب سایت های کتابخانه های ملی کشورهای اسلامی در منطقه خاورمیانه می پردازد کاملا مطابقت و همسویی دارد.

۱۸۰ مدیریت منابع اطلاعاتی وب

جدول ۱. سیاهه ارزیابی کاربردپذیری وب سایت نهاد کتابخانه‌های عمومی کشور از نظر میزان رعایت مؤلفه‌ها (به تفکیک)

نهاد کتابخانه‌های عمومی کشور		اعتبار اطلاعات	
ردیف	تواتر	مؤلفه‌ها	
۱	۱,۴۵,۱۰,۱۵,۱۷,۱۸,۱۹,۲۳,۲۵	بیان نام مولف	-
۲	۱,۱۰,۳۵,۱۵,۱۶,۱۷,۱۸,۱۹,۲۳	بیان ویژگی‌های مولف (صلاحیت, شهرت, اعتبار و...)	-
۳	۶,۸,۱۰,۴۵,۱۹,۴۲, ۱	چگونگی تماس با سازمان یا شخص مسئول محتویات سایت	✓
۴	۱,۶,۱۵,۱۸,۲۱,۲۳,۲۶	دارا بودن تأییدیه رسمی از سوی سازمان یا فرد مسئول محتوی سایت	-
۵	۶,۲۵,۸,۱۵,۱۷,۴۲ و ۱	ذکر مشخصات سازمان یا شخص مسئول محتویات سایت	✓
۶	۱,۸,۱۵,۲۵,۳۷,۴۲	ارائه فهرستی از کارگزاران اصلی و مشخصات آن‌ها	-
۷	۲۴,۲۵,۲۹,۴۰, ۲۷,۱,۸,۱۰, ۱۵, ۱۹	ذکر راهی برای تماس با مولف صفحه (نشانی پستی, شماره تلفن و پست الکترونیک و ...)	✓
۸	۱,۱۵,۱۶,۱۸,۱۹,۲۰	بیان راهی برای اثبات ویژگی‌های مولف (تجارب وی در زمینه موضوعی خاص, عضویت وی در سازمان‌های حرفه‌ای و...)	-
۹	۱,۸,۲۷,۱۵,۱۷,۱۹	ذکر نام سازمان مسئول محتویات سایت	✓
۱۰	۶,۱۸,۱۹, ۱	بیان ماهیت حامی	-
۱۱	۷,۴۲, ۱,۸	ذکر مدت زمان تاسیس سازمان	-
۱۲	۱,۸,۱۵,۳۵,۲۲	بیان حق مالکیت معنوی منابع اطلاعاتی عرضه شده در سایت (ذکر نام)	✓
۱۳	۱,۸,۱۵,۴۲	بیان نام سازمان پشت پرده سایت	-
۱۴	۱,۸,۱۰, ۱۵	ارائه سیاهه‌ای از نام سازمان‌ها یا سایت‌هایی که این سایت را توصیه می‌کنند	-
۱۵	۱,۱۵,۲۹,۱۶,۲۵	مشخص بودن هدف از طراحی و نشر صفحه	-
۱۶	۱,۸,۱۰, ۱۵	ارائه فهرستی از اسامی و مشخصات افراد مسئول نظارت بر سازمان	-
۱۷	۱, ۱۵	ارائه فهرستی از منابع چاپی منتشر شده توسط سازمان	-
۱۸	۲۷,۲۲,۴۰	قرار دادن نام و آرم سازمان سایت در هر صفحه	✓
۱۹	۳۰	نشان دادن آرم سازمان بصورت واضح و برجسته	✓
۲۰	۲۲	قرار دادن نام سازمان سایت در بالای صفحه	✓
۲۱	۱, ۱۵	ذکر نام دارنده حق مالکیت معنوی	✓
۲۲	۹	قراردادن پرچم کشور	✓
۴۵	جمع امتیازها	۱۰۷	
نهاد کتابخانه‌های عمومی کشور		صحت اطلاعات	
ردیف	تواتر	مؤلفه‌ها	
۲۳	۸,۲۵,۳۶,۱۶,۱۹, ۲۷,۲۲,۴۰, ۱,۶, ۱۰	نبود خطاهای املائی, گرامری و تاپی	✓
۲۴	۱,۱۰,۲۲,۱۵,۱۷,۲۴,۱۸,۲۳	ذکر نام و مشخصات کتابشناختی منبع اصلی	-
۲۵	۱,۱۰,۱۵,۱۹,۲۸,۱۷, ۱۹	وجود شاخصی دال بر بررسی صحت و سقم اطلاعات توسط ویراستار یا مصحح در طول فرآیند بازنگری مجدد منابع	-
۲۶	۸,۳۵, ۲۷,۳۹, ۱, ۱۰	عنوان بندی روشن و واضح گراف, نمودار یا جداول موجود در صفحه	✓
۱۷	جمع امتیازها	۳۲	

جدول 1 سیاهه ارزیابی کاربردپذیری وب سایت نهاد کتابخانه‌های عمومی کشور از نظر میزان رعایت مؤلفه‌ها (به تفکیک)

ارزیابی کاربردپذیری وبگاه نهاد کتابخانه‌های عمومی کشور ۱۸۱

نهاد کتابخانه‌های عمومی کشور		روزآمد بودن اطلاعات	
ردیف	تواتر	مؤلفه‌ها	
۲۷	۷,۱۰,۱۵,۱۷,۱۸,۱۹,۲۳,۲۵,۱	ذکر تاریخ آخرین تجدید نظر در محتویات صفحه	-
۲۸	۸,۱۶,۱۷,۱۸,۱۹,۲۳,۱,۱۰,۳۵	ذکر نخستین تاریخ قرار گرفتن منبع اطلاعاتی (با هر فرمتی) بر روی صفحه وب	-
۲۹	۶,۸,۱۰,۱۵,۱۹,۴۲,۱	ذکر فواصل زمانی به روز کردن اطلاعات دارای حساسیت زمانی	-
۳۰	۱,۶,۱۵,۱۸,۲۱,۲۳,۲۶	وجود اطلاعات آماری در صفحه	✓
۳۱	۱,۶,۲۵,۸,۱۵,۱۷,۴۲	بیان تاریخ گردآوری آمار	-
۳۲	۱,۸,۱۵,۲۵,۳۷,۴۲	ارائه تاریخ‌ها در یک فرمت بین‌المللی	-
۷	۲۶	جمع امتیازها	
نهاد کتابخانه‌های عمومی کشور		سطح پوشش مخاطبان خاص	
ردیف	تواتر	مؤلفه‌ها	
۳۳	۴۴,۴۰,۱,۶,۱۰,۸,۱۵,۱۶,۱۸,۲۸	مشخص بودن نام و نوع منابع موجود در صفحه	✓
۳۴	۱,۶,۱۰,۸,۱۵,۱۸,۲۸	تعیین مخاطبان خاص صفحه	-
۳۵	۱,۶,۱۰,۱۵	درج زمان تخمینی تکمیل صفحه در دست ساخت	-
۱۰	۲۱	جمع امتیازها	
نمایهای تعاملی و تبادل		نهاد کتابخانه‌های عمومی کشور	
ردیف	تواتر	مؤلفه‌ها	
۳۶	۲۷,۳۰,۱,۸,۴,۱۵,۱۱,۲۶,۳۸	وجود نظام مشخصی جهت بازخورد کاربران (feedback)	✓
۳۷	۱,۸,۴,۱۵,۲۶,۳۸	وجود نظام مشخصی برای کاربران برای درخواست اطلاعات بیشتر از سازمان	-
۳۸	۱,۲۷,۱۵,۲۶,۳۸	ذکر مشخصه زمانی لازم برای دریافت پاسخ از سوی سازمان	-
۳۹	۱,۲۶,۳۸	وجود امکان عضویت در سایت	-
۴۰	۱,۲۶,۳۸	وجود نظام مشخصی جهت عضویت کاربران در سایت	✓
۴۱	۱,۱۵	آگاهی کاربر از وجود مکانیزم کوکیز در سایت	-
۴۲	۲۷	پاسخگویی سریع و دقیق سیستم FAQ به پرسشهای کاربر	-
۴۳	۲۷	قابلیت پاسخگویی به سئوالات کاربران	-
۴۴	۲۷	آگاهی کاربران از زمان دریافت پاسخ	-
۱۲	۳۱	جمع امتیازها	
عینیت اطلاعات		نهاد کتابخانه‌های عمومی کشور	
ردیف	تواتر	مؤلفه‌ها	
۴۵	۱,۱۰,۱۵,۱۶,۲۳,۴	مشخص بودن ارتباط میان شخص مولف یا سازمان و فرد مسئول محتویات سایت	-
۴۶	۱,۱۵,۱۶,۴,۲۳	مشخص بودن دیدگاه مولف	-
۴۷	۱,۸,۱۵,۴,۳۸	مشخص بودن دیدگاه شخص یا سازمان مسئول تهیه اطلاعات	✓
۴۸	۱,۱۰,۱۵,۱۸	نبود در تبلیغات در صفحه	✓
۴۹	۱,۶,۱۷,۱۵,۲۷,۳۰	ارائه مطالبی در خصوص اهداف شخص یا سازمان گردآورنده اطلاعات	-
۵۰	۱,۱۵	توجه اطلاعات غیر مرتبط با خدمات سازمان در صفحه	-

✓	وجود وجه تمایزی میان محتوی اطلاعاتی و محتوی تقریباتی	۱۵,۳۹,۲۷	۵۱
✓	شبه نبودن اطلاعات به تبلیغات از لحاظ نوع طراحی	۲۷	۵۲
✓	اجتناب از بکاربردن مفاهیم و اطلاعات نا مرتبط	۴۰,۳۹	۵۳
۱۵	۳۴	جمع امتیازها	
نهاد کتابخانه‌های عمومی کشور		ناوبری	
	الف: ویژگیهای عنوان مرورگر	تواتر	ردیف
✓	مبین نام سازمان یا فرد مسئول محتویات سایت	۱,۱,۲۱	۵۴
✓	مبین صفحه اصلی بودن صفحه	۱,۱۵,۲۶	۵۵
✓	کوتاه بودن عنوان مرورگر	۱,۱۵	۵۶
✓	منحصر بودن عنوان مرورگر برای سایت	۱۵	۵۷
۹	۹	جمع امتیازها	
	ب: ویژگیهای عنوان صفحه	تواتر	ردیف
✓	مبین متعلق بودن صفحه اصلی بودن صفحه	۱,۱۰,۳۴,۱۵,۱۶,۲۱	۵۸
-	مبین متعلق بودن صفحه به یک سایت مشخص (حداقل توسط یک آرم)	۱,۱۰,۲۷,۱۵,۲۶	۵۹
✓	کوتاه بودن عنوان صفحه	۷۰۱	۶۰
✓	منحصر بودن عنوان صفحه برای سایت	۳۴,۴۰,۴۴,۱۵	۶۱
✓	بکاربردن کلمات کلیدی و مهم برای عناوین	۲۲,۲۷	۶۲
✓	بکارگیری عنوان مناسب جهت بیان دقیق محتوی متن	۲۷,۲۹,۲۲	۶۳
۱۷	۲۲	جمع امتیازها	
	ج: پیوندهای متنی و فرا متنی	تواتر	ردیف
✓	پکنواختی ظاهر پیوندهای تکراری	۱,۶,۱۰,۴,۱۵,۳۸,۲۵,۲۶	۶۴
✓	امکان ناوبری آسان میان صفحات	۲۷,۳۰,۲۴,۱۶,۱۰,۴,۱۵,۲۱,۲۵	۶۵
✓	قرارگیری هماهنگ پیوندهای مستقیم داخلی در صفحه	۱,۶,۳۰,۱۰,۴,۱۵,۲۱,۲۶	۶۶
✓	وجود تناسب میان عنوان پیوند با آنچه پیوندبان ختم می شود	۱,۱۵,۱۷,۲۱,۲۵,۲۶,۳۵	۶۷
-	ادامه بکار پیوندها متناسب با اهداف از پیش تعریف شده	۱,۱۵,۱۶,۱۷,۱۸,۲۶	۶۸
✓	قابل درک بودن ساختار سایت / صفحه برای کاربران	۲۷,۴۰,۳۳,۱۶,۴,۱۵,۲۱,۲۶	۶۹
✓	تناسب چیدمان صفحات مختلف سایت با یکدیگر	۲۷,۱۰,۱۵,۲۱,۲۵	۷۰
-	وجود مابتهایی در صفحه اصلی جهت دسترسی هرچه آسانتر و موثر تر کاربران به صفحات پرمخاطب سایت	۱,۱۵,۲۱,۲۶,۳۸	۷۱
✓	وجود امکان دسترسی از صفحه اصلی سایت به بخش های اصلی آن	۱,۱۰,۳۹,۱۵,۲۱,۲۶	۷۲
✓	اجتناب از بکاربردن متن های زیر خط دار در کنار پیوندها	۱,۱۰,۱۵,۲۱	۷۳
-	عنوان دهی مناسب بوک مارک ها	۱۵,۲۱,۲۵	۷۴
-	امکان انتخاب اقلام اطلاعاتی مورد نظر از روی فهرست به جای تایپ آنها	۱,۱۵,۱۹	۷۵
✓	قرار دادن پیوند آرم سایت در صفحه اصلی	۱۵,۲۶	۷۶
-	شناساندن پیوندها از طریق زیر خط دار کردن آنها یا کاربرد نوعی رنگ خاص	۳۰	۷۷
✓	وجود تناسب بین عنوان صفحه و پیوند	۲۷	۷۸
-	نمایان بودن پیوندهای بازدید شده و بازدید نشده با استفاده از تغییر رنگ پیوند	۴۰	۷۹

ارزیابی کاربردپذیری وبگاه نهاد کتابخانه‌های عمومی کشور ۱۸۳

✓	قرار دادن پیوند صفحه اصلی جهت مشخص بودن صفحه اصلی	۲۷,۴۴	۸۰
-	منطقی بودن تعداد پیوندها	۴۴	۸۱
-	پیوند دادن تصویر به صفحه مربوطه	۳۰	۸۲
✓	امکان شناسایی آسان پیوندها	۲۴,۲۷	۸۳
✓	استفاده از عبارات مناسب برای پیوندها (عدم استفاده از عبارت , klik here , more)	۲۴	۸۴
-	قرار دادن پیوند متنی در ابتدای پاراگراف	۲۲,۳۰	۸۵
✓	وجود تناسب بین عنوان صفحه و پیوند	۲۴	۸۶
-	استفاده از علائم دیداری مانند رنگ، اندازه برای نشان دادن ارتباط میان پیوندها	۲۲	۸۷
-	بکارگیری متن کافی برای توضیح پیوند	۴۰	۸۸
۶۸	۹۲	جمع امتیازها	
	۵: نشانه اینترنتی صفحه	تواتر	ردیف
-	درج نشانه اینترنتی صفحه در بدنه اصلی آن	۱,۱۵,۲۱	۸۹
✓	مختصر بودن نشانه اینترنتی صفحه	۲۷	۹۰
✓	عدم تغییر در نشانه اینترنتی سایت	۲۷	۹۱
✓	کاربر پسند بودن نشانه اینترنتی صفحه	۳۰	۹۲
۳	۶	جمع امتیازها	
	۵: نقشه سایت یا نمایه	تواتر	ردیف
-	وجود نقشه سایت یا نمایه ها بر روی صفحه اصلی و صفحات پیوندی	۲۷,۳۹,۱,۶,۱۵,۲۱,۲۶,۳۷	۹۳
-	دارا بودن موضوعات اصلی سایت	۱,۶,۱۵,۲۱,۲۶,۳۷	۹۴
-	امکان خواندن آسان نمایه ها یا نقشه سایت	۱,۶,۱۵,۲۱,۳۷	۹۵
-	سازماندهی منطقی نمایه ها یا نقشه سایت	۱,۲۷,۶,۱۵,۲۱,۳۷	۹۶
۰	۲۵	جمع امتیازها	
	ی: موتور جستجوی داخلی	تواتر	ردیف
✓	وجود بک موتور جستجوی داخلی	۱,۴,۱۵,۲۱,۲۵,۲۸	۹۷
-	مرتبط بودن اطلاعات بازیابی شده توسط موتور جستجوی داخلی	۱,۴,۱۵,۲۱,۲۵,۲۸	۹۸
۶	۱۲	جمع امتیازها	
نهاد کتابخانه های عمومی		نماهای غیرمتنی	
	مؤلفه ها	تواتر	ردیف
✓	عدم استفاده از تصاویر متحرک (انیمیشن) بی مورد	۲۷,۳۹,۱,۶,۱۵,۲۱,۲۶,۳۵,۴۲,۳۷	۹۹
-	بکارگیری تصاویر گرافیکی، فایل های صوتی و تصویری به منظور افزایش کارایی سایت	۲۲,۴۰,۱,۱۰,۱۵,۲۱,۲۶,۳۷,۳۵	۱۰۰
-	بیان نام نرم افزار خاص مورد نیاز و چگونگی دسترسی به آن	۱,۱۰,۱۵,۲۱,۲۶,۳۷,۳۵	۱۰۱
-	وجود جایگزینی برای فایل های نیازمند نرم افزار خاص جهت دسترس پذیری اطلاعات برای کلیه کاربران	۱,۱۰,۱۵,۲۱,۲۶,۳۷,۳۵	۱۰۲
-	بیان نام نرم افزار مرورگر مورد نیاز، یا ویرایش خاصی از آن، جهت دسترسی به	۱,۱۰,۱۵,۲۶,۳۵,۳۷	۱۰۳

۱۸۴ مدیریت منابع اطلاعاتی وب

صفحات وب (در صورت لزوم)			
-	وجود جایگزین متنی برای تصاویر و گرافیک های موجود در صفحه برای کلیه کاربران	۱,۱۰,۱۵,۲۷,۲۱,۲۶,۳۷,۳۵	۱۰۴
✓	عدم استفاده از تکنیک فلش (علامت چشمک زن)	۱,۱۰,۳۰,۲۱,۳۵,۳۷,۴۲	۱۰۵
-	آگاهی کاربر از بار شدن فایلی حجیم در صورت پیروی از یک پیوند	۱,۱۵,۲۱,۲۶,۳۵,۳۷	۱۰۶
✓	مرتبط بودن تصاویر گرافیکی با محتوی متن	۲۷	۱۰۷
-	اضافه کردن متون به تصاویر برای درک بیشتر	۲۷	۱۰۸
-	استفاده از تگ ALT مناسب برای تصاویر	۳۰	۱۰۹
-	هدفمند بودن تصاویر گرافیکی	۳۲	۱۱۰
۱۸	۶۴	جمع امتیازها	
تعداد کتابخانه های عمومی کشور		دسترس پذیری	
	مؤلفه ها	تواتر	ردیف
✓	دسترس پذیر بودن با استفاده از مرورگر اینترنت اکسپلورر (Internet explorer 6.0)	۱,۶,۱۰,۱۵,۱۶,۲۵,۴۱,۳۷,۴۳,۳۵	۱۱۱
✓	دسترس پذیر بودن با استفاده از مرورگر نت اسکپ نیویگیتور (Netscape Navigator 6.5)	۱,۶,۱۰,۱۵,۱۶,۲۶,۴۱,۳۷,۴۳,۳۵	۱۱۲
-	دسترس پذیر بودن با استفاده از مرورگر موزیلا برد (Mozilla Bird 0.7)	۳۷,۴۳,۳۵, ۱,۶,۱۰,۱۵,۱۶,۲۶,۴۱	۱۱۳
-	دسترس پذیر بودن با استفاده از مرورگر اپرا (Opera 7.2)	۳۷,۴۳,۳۵, ۱,۶,۱۰,۱۵,۱۶,۲۶,۴۱	۱۱۴
✓	دسترس پذیری سایت از طریق موتورهای جستجوی عمومی	۱,۱۵,۱۶,۱۷,۴۳,۲۵,۲۶	۱۱۵
-	اندازه صفحه کمتر از ۵۰ کیلو بایت	۱,۴,۳۹,۲۶,۳۵,۳۸,۴۱	۱۱۶
✓	امکان پشتیبانی از سکوی عملیاتی کاربر (وبندوز) me,۹۸,۲۰۰۰, لینوکس (REDHAT)	۱۰,۶,۱۵,۲۶	۱۱۷
✓	استفاده از قلم (فونت) های استاندارد	۱۲, ۴, ۱۵, ۳۵, ۳۸, ۲۷, ۳۰, ۲۴, ۴۴, ۳۲, ۳۴	۱۱۸
✓	قابل رویت بودن تمامی اجزای سایت	۱۵,۳۸	۱۱۹
-	دسترسی کاربر به منابع مورد نیاز در کمتر از ۳ بار کلیک کردن	۲۷	۱۲۰
-	شبه بودن اطلاعات به هرم واژگون (دستیابی سریع به اطلاعات مهمتر)	۲۷	۱۲۱
✓	سهولت استفاده از سایت برای همه کاربران (مبتدی و متخصص)	۲۷ و ۲۴	۱۲۲
✓	عدم استفاده از تصاویر متحرک (انیمیشن) بی مورد	۳۹ و ۳۰	۱۲۳
۴۸	۷۷	جمع امتیازها	
تعداد کتابخانه های عمومی		کارآمدی	
	مؤلفه ها	تواتر	ردیف
✓	استفاده از رنگ های استاندارد در طراحی صفحات	۴۲,۳۷,۳۵ و ۱۶ و ۲۱ و ۱۰ و ۳۴ و ۳۲ و ۳۹ و ۲۴ و ۲۷	۱۲۴
✓	عنوان بندی مناسب قابلیت های سایت / صفحه	۲۵,۲۶, ۱,۱۵,۲۱,۱۷	۱۲۵
✓	امکان تشخیص سایت بر اساس عنوان دامنه صفحه (Domin)	۱۵,۲۸,۲۵,۱۹,۳۵	۱۲۶
✓	تناسب زبان سایت با فرهنگ و روحیات کاربر	۱,۱۵,۳۹,۲۶,۲۸,۳۸	۱۲۷
-	وجود تعاریف حاشیه ای برای تشریح اقلام اطلاعاتی موجود بر روی صفحه	۱,۱۵,۲۶,۳۸	۱۲۸

ارزیابی کاربردپذیری وبگاه نهاد کتابخانه‌های عمومی کشور ۱۸۵

-	اطلاع کاربر از عملیات در حال انجام	۳۰,۲۲,۱,۱۵,۲۶,۳۸	۱۲۹
-	تناسب اطلاعات موجود در صفحات یا ماموریت سایت	۱,۱۵,۲۶,۳۸	۱۳۰
✓	امکان انجام تمامی قابلیت های سایت در درون سایت	۱۵,۲۶,۳۷	۱۳۱
-	وجود قابلیت لغو عملیات صورت گرفته در سایت	۱۰,۱۵,۲۵	۱۳۲
-	وضوح موقعیت کاربر در سایت	۱,۱۵,۲۶	۱۳۳
✓	در نظر گرفتن پیشینه ذهنی کاربر	۱,۳۷,۳۹,۲۶	۱۳۴
-	جذابیت سایت	۱۷,۲۶	۱۳۵
-	در نظر گرفتن ابزارهایی جهت کمک به کاربر	۱,۱۵,۲۶,۳۸	۱۳۶
-	ارائه عنوان کامل تمام سرواژه های مهم بکاررفته در صفحات	۱۵,۳۴,۲۵,۳۰	۱۳۷
✓	طومارنوردی آسان صفحه خانگی	۶,۴۲	۱۳۸
✓	امکان چاپ اطلاعات بدون اعمال تغییرات در تنظیم سیستم رایانه	۲۵,۱۵	۱۳۹
✓	بیان تعداد مراجعان سایت در یک بازه زمانی	۱۵	۱۴۰
-	استفاده از عبارات کوتاه برای تشریح موارد موجود در صفحه	۲۷,۴۴,۲۶	۱۴۱
-	ایجاد جلب توجه کاربران از طریق نوع طراحی (نوع رنگ و فونت و ...)	۲۷,۳۴	۱۴۲
✓	بهم فشرده نبودن مطالب در یک صفحه و ایجاد فضای خالی بین آنها	۲۷,۳۹,۴۰	۱۴۳
-	دسترسی آسان به بخش Help	۳۹	۱۴۴
-	ارائه قیمت پیشنهادی به کاربر در صورت انتفاعی بودن سایت	۳	۱۴۵
-	بالا بودن سرعت بارگذاری	۳۰,۲۷	۱۴۶
✓	قابلیت استفاده از عملکرد جستجو	۲۷,۳۹,۳۴,۲۲	۱۴۷
-	امکان ناوبری توسط توروک	۲۹	۱۴۸
-	امکان ناوبری توسط جستجو	۲۹,۲۷	۱۴۹
-	دارا بودن خطی متنی اختصاصی	۲۷,۴۴	۱۵۰
-	قرار دادن مطالب در کمتر از ۳ صفحه	۲۴	۱۵۱
✓	رعایت سلسله مراتب مطالب در صفحات (بر حسب تاریخ و عنوان و ...)	۲۹	۱۵۲
✓	بکار بردن عدد برای نشان دادن ارقام (۲ به جای دو)	۴۴	۱۵۳
✓	عدم استفاده از فایل های pdf بدلیل کاهش جذابیت (به جز برای اسناد)	۴۴	۱۵۴
-	قابل رویت بودن پیام های اخطار	۳۹	۱۵۵
-	پایین بودن زمان داتلود	۳۴,۴۴	۱۵۶
۵۱	۹۹	جمع امتیازها	
نهاد کتابخانه‌های عمومی کشور		ویژگی‌های ظاهری	
	مؤلفه ها	تواتر	ردیف
✓	صفحه آرای مناسب	۲۲,۳۴,۴۰	۱۵۷
✓	وجود هماهنگی بین طرح ها و صفحات	۲۷	۱۵۸
✓	وجود هماهنگی بین رنگ ها و سبک نگارش	۳	۱۵۹
✓	بکارگیری حاشیه در اطراف متن	۲۴	۱۶۰
۶	۶	جمع امتیازها	

در مجموع نتایج به دست آمده از ارزیابی کاربرد پذیری وب سایت نهاد کتابخانه های عمومی کشور به تفکیک معیارهای 11 گانه ی موجود در این پژوهش در مقایسه با وب سایت شاخص فرضی با 663 امتیاز، در سطحی معادل 52 درصد (346 امتیاز)، قرار دارد و این امر گر چه بیان گر آن است که تنها نیمی از استانداردها در وب سایت نهاد کتابخانه های عمومی کشور رعایت گردیده است و این امر با وب سایت شاخص فاصله زیادی دارد لذا لازم است در طراحی این وب سایت دقت و توجه بیش تری به عمل آید افزون بر این با عنایت به کاربردی بودن این پژوهش پژوهش گران در صدند تا با انعکاس نتایج پژوهش حاضر، کتابداران و طراحان وب سایت ها را از نقاط ضعف و قوت وب سایت مطلع ساخته تا با به کارگیری عناصر لازم و ویژگی های مناسب در طراحی وب سایت فضایی را ایجاد نمایند که کاربر به محض ورود به وب سایت امکان بهره گیری از تمامی اطلاعات و خدمات مرتبط با این، محمل بدون وجود هر گونه محدودیتی را داشته باشد.

6- پیشنهاد های پژوهش

با توجه به یافته های پژوهش در باب رعایت معیارهای مطرح در کاربرد پذیری هر چه بیشتر و بهتر وب سایت نهاد کتابخانه های عمومی کشور توجه به رعایت معیارها و مولفه های: بیان نام مولف و ویژگی های او (صلاحیت شهرت و اعتبار و...)، ذکر فواصل زمانی به روز کردن اطلاعات دارای حساسیت زمانی درج زمان تخمینی تکمیل صفحه در دست ساخت نماهای تعاملی و تبدالی، وجود امکان عضویت در سایت، مشخص بودن ارتباط میان شخص مولف با سازمان و فرد مسئول محتویات سایت منحصر بودن عنوان صفحه برای سایت نمایان بودن پیوندهای بازدید شده و بازدید نشده با استفاده از تغییر رنگ پیوند درج نشانه اینترنتی صفحه در بدنه اصلی آن بیان نام نرم افزار خاص مورد نیاز و چگونگی دسترسی به آن بیان تعداد مراجعان سایت در یک بازه زمانی توصیه می گردد

منابع

1. اصغری، پوده احمد رضا 1380. بررسی عناصر و ویژگی های مطرح در طراحی وب سایت کتابخانه های دانشگاهی. پایان نامه کارشناسی ارشد کتابداری و اطلاع رسانی دانشکده علوم تربیتی و روانشناسی، دانشگاه فردوسی مشهد
2. الکساندر ژانت 1383. شناخت وب: چگونه اطلاعات موجود بر روی وب را ارزیابی نموده و صفحاتی اینگونه را پدید آوریم، ترجمه صدیقه محمد اسماعیل، تهران: دبیزش.
3. جانسون، استیو 1382. صفحات سفارش و دریافت مواد منابع کتابخانه ای بر روی وب. ترجمه صدیقه محمد اسماعیل اطلاع شناسی، 1 (پاییز): 187-199.
4. حاجی زین العابدینی، محسن مکتبی، فرد، لیلا عصاره، فریده 1384 تحلیل پیوندهای وب سایت های کتابخانه های ملی جهان مجله مطالعات تربیتی و روانشناسی دوره 7، شماره 1، ص 193-173

5. خوانساری، جیران ساختن سایت های وب موفق برای کتابخانه های کوچک دانشگاهی. صنعت برق 54 (آبان) ص. 24-28.
6. دایره المعارف کتابداری و اطلاع رسانی ج. 1 و 2. 1381 تهران: نهاد کتابخانه های عمومی کشور، 1381.
7. رضایی شریف آبادی، سعید، فرودی، نوشین 1381. ارزیابی صفحات وب کتابخانه های دانشگاهی ایران و ارائه الگوی پیشنهادی فصلنامه کتاب دوره سیزدهم 4 (زمستان) ص. 12-19.
8. عصاره، فریده، مرادمند، علی 1384. شناسایی ویژگی های عمده در طراحی وب سایت های کتابخانه های ملی جهان به منظور ارائه الگویی مناسب جهت ارتقاء کیفی وب سایت نهاد کتابخانه های عمومی کشور. فصلنامه اطلاع شناسی پاییز و زمستان شماره 1 و 2. ص. 170-190
9. عصاره فریده 1381 معیارهای ارزیابی منابع اینترنتی فصلنامه کتاب، دوره سیزدهم 2 (تابستان). ص. 61-73
10. محمد اسماعیل صدیقه 1383. بررسی کاربرد پذیری وب سایتهای دانشگاه های صنعتی کشور. پایان نامه دکتری کتابداری و اطلاع رسانی تهران دانشگاه آزاد اسلامی واحد علوم و تحقیقات
11. محمد اسماعیل صدیقه 1384 بررسی کاربرد پذیری وب سایت های دانشگاه های صنعتی کشور فصل نامه کتاب دوره 16 شماره 1.
12. نویدی، فاطمه 1386. ارزیابی دسترس پذیری وب سایت وزارتخانه های دولت جمهوری اسلامی. ایران پایان نامه کارشناسی ارشد علوم کتابداری و اطلاع رسانی دانشگاه تربیت مدرس
13. Alexander, J. E. ,Tate, M. A.1999. Web Wisdom: who to evaluate and create information quality on the . (web. London, Mahwah, New Jersey. LEA: (Lawrence Erlbaum Associates
14. Back, S. E.1997. Evaluation criteria: the good, The Bad and The Ugly: Or, Why It's a Good Idea to . "Evaluate Web Sources". [on-line] Available: <http://lib.nmsu.edu/instruction/evalcrit.html>
15. Barker. J. 2004 Finding Information on the Internet: A Tutorial University of California. [on-line]. Available: <http://www.lib.berkeley.edu/TeachingLib/Guides/Internet/Evaluate.html>
16. Barker. J. 2003. Evaluating Web Pages: Techniques to Apply Questions to Ask." VC Berkeley - . Teaching Library Internet Workshops, 12sep.2003.[on-line]. Available: <http://www.Lib.berkeley.edu/Teaching Lib/Guides/Internet/Evaluated.html>
18. Berger. P. 1999 Web Evaluation Guide: Tramline, Incorporated. [on-line]. Available: .

<http://www.infosearcher.com/cybertours/tours/touro4/-tourlaunch1.html>

Bertot J.C. et al. 2006 Functionality, usability and accessibility: Interactive user-centered evaluation . 19 strategies for digital libraries”, Performance Management and Metrics, Vol. 7

ص: 187

- Braynik. G. 2003. Atomic Web Usability Evaluation: What Need to be done?. (29jan.2003): 1–16. .20
[on–line]. Available: <http://Usable.binghamton.edu.Atomic Thesis. html>
- Clyde. L. A. 1996. The Library as information providers, The Homepage The Electronic Library. Vol. .21
.14 no.6 pp:549–558
- Engle. M. 1996. Evaluating Websites: Criteria and Tools. New York Library Association Conference, .22
Saratoga Springs, Ny. (October 1996): 1–3.[on–line]. Available: <http://www.Library.Cornel.edu/Okuref/research/Webeval.html>
- .The Essential Web Site Usability Checklist. 2007. [on–line]. Available: www.dailybits.com .23
- Hupp.J.2008. Test Your Web Site: A 57– Point Checklist for Maximum Usability .[on– online]. . 24
.Available: www.virtualhosting.com
- .Jafari. A. Optimizing Campus Web Sites. EDUCASE QUARTERLY, No.2 2000:56–58 .25
- .Leggett. D. Quick Usability Checklist.2009, [on–line]. Available: www.uxbooth.com .26
- Meyers. P. J. 25–Point Web Site Usability Checklist, 2008. [on–line]. Available: www.usereffect.com .27
- Nadler.D.M. and Furman.V.M. (2001). Access board issues final standards for disabled access under .28
.Section 508 of Rehabilitation Act, Government Contract Litigation Reporter, Vol. 14 No. 19, p. 14
- National Science foundation (NSF).Universal Design of College Algebra, 2008. [on– line]. Available: .29
www.usablealgebra.landmark.edu
- Nielsen.J.Coyne, K, Tahir, Marie, (2001). Make It Usable–Web Site Usability Magazine, (6 feb.2001): .30
.1–5.[on–line]. Available: <http://www.Pcmag.com/article2/0,4149,33821,00. asp>, 23jan.2003
- Osareh.F. (2003). "A notice on content of library information science (LIS) schools websites, Libri (3): .31
262–265
- .The SEO File: Usability Checklist 2009. [on–line]. Available: www.theseofiles.co.uk .32
- Sloan, John. Rampo College Web Design Standards. Rampo College of New jersey Statements and .33
policies.2001.[on–line].Available: <http://www.guide.rampo.edu/1/content/Webstandards.html>

,Stueart. R.D. Library and Information Center management fifth edition. Endewood .34

ص: 188

- Sullivan, T. Matson R. Barriers to use: Usability and Content Accessibility on the Web's Most Popular .35
.Sites. Interdisciplinary PhD. program in information Science, University of North Texas, Denton, 2000
- Tungar.M. N. Heuristic Evaluation.2002.[on-line]. Available: <http://www.cc.gatech.edu/~manas/cs8803a/Heuristic.pdf> . 36
- .Usability Guidelines Information Science Technology, 2004. [on-line]. Available: www.web.mit.edu .37
- .Web Usability Checklist. The College of New Jersey. 2008. [on-line]. Available: www.teng.edu" .38
- Weibel .S.L. The World Wide Web and Emerging Internet Resource Discovery Standards for Scholarly .39
.Litrature. Library Trends.vol.43, No.4 1995:627-664
- .Wilson, S. 1995 World Wide Web Design Guide. Indianapolis. IN: Hayden Books, 1995 .40
- Zaphiris. P. Darin, Ellis, R. Website Usability and Content Accessibility of The Top USA Universities .41
2001. Dertoit, MI: Institute of Geronotology and Dept of Industrial and Manufacturing Engineering Wayne
.State University, 2001
- Web Usability Tips to Attract and Retain Web Visitors 2009,:[on-line].Available: 50 . 42
www.doshdosh.com

هدف: هدف کلی پژوهش دستیابی به فرایند پردازش و سازماندهی وبگاه‌ها با استفاده از قواعد و استانداردهای مورد استفاده در سازمان اسناد و کتابخانه ملی ایران است.

جامعه آماری: جامعه آماری این پژوهش را 50 وبگاه تشکیل می‌دهند که در 20 رده موضوعی از وبگاه پارس ایندکس و به صورت تصادفی انتخاب شده‌اند.

روش شناسی: این پژوهش به دوروش کتابخانه‌ای و پیمایش توصیفی انجام شده است- جهت گردآوری اطلاعات از سیاهه‌های واری و برای تجزیه و تحلیل داده‌ها از آمار توصیفی استفاده شده است.

نتایج: نتایج نشان می‌دهند که استاندارد سازی موضوع سازماندهی و پردازش وب‌گاه‌ها به عنوان نوعی از منابع الکترونیکی و بومی سازی آن در سازمان اسناد و کتابخانه ملی ایران امکان پذیر است.

کلیدواژه‌ها: وب‌گاه‌ها سازماندهی امکان‌سنجی سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران

دکتر رضا خانی پور (1) | محبوبه قربانی (2) | سهیلا فعال (3)

مقدمه

امروزه شبکه جهانی اینترنت منابع اطلاعاتی متنوعی را در دسترس قرار داده است. وب هر روزه پهنه بیش تری از وسعت دنیا را در بر گرفته و در دور افتاده ترین نقاط جهان نیز رسوخ نموده است. شبکه جهانی وب با دسترسی پذیری خود امکانات زیادی را برای مخاطبان فراهم نموده است. دسترسی به آخرین اخبار دنیا و اطلاعات کشفیات و نوآوری ها آشنایی با فرهنگ و تمدن ملل مختلف، دریافت اطلاعات مورد نیاز روزانه مانند آب و هوا و اخبار وقایع مهم از امتیازات دسترسی به وب است. این گستردگی و دسترسی پذیری در کنار محاسنی که دارد کاربران را با مشکلاتی روبه رو ساخته است. اطلاعات در لابلای صفحات وب نهفته و نیاز به جستجوهای دقیقی برای کشف آن ها وجود دارد.

عمده ترین مشکل کاربران، وب دست یابی به بهترین و با کیفیت ترین اطلاعات مورد نیاز و حصول جامعیت و مانعیت در بازبایی اطلاعات در زمینه های تخصصی است (علی محمدی 1381)

ص: 191

1- دکترای کتابداری و اطلاع رسانی عضو هیات علمی و مدیر کل پردازش و سازماندهی سازمان اسناد و کتابخانه ملی ایران -R

Khanipour@nlai.ir

2- دانشجوی دکترای کتابداری و اطلاع رسانی معاون اداره کل پردازش و سازماندهی سازمان اسناد و کتابخانه ملی ایران -m

ghorbani@nlai.ir

3- کارشناس ارشد کتابداری و اطلاع رسانی رئیس گروه سازماندهی منابع غیر کتابی سازمان اسناد و کتابخانه ملی ایران -s

faal@nlai.ir

کتابخانه ها و مراکز اطلاع رسانی به عنوان متولیان گردآوری سازماندهی و اشاعه اطلاعات، خود را موظف به گردآوری سازماندهی اطلاعات موجود بر روی وب می دانند یکی از چالش های فراروی کتابداران نحوه سازماندهی منابع وب بوده است تلاش های گسترده ای در خصوص سازماندهی اطلاعات بر روی وب صورت گرفته است تهیه انواع استانداردها و دست نامه ها و ابداع فراداده های متنوع سازماندهی منابع وب از جمله تلاش های صورت گرفته در جهت سازماندهی منابع وب است. (نشاط 1382)

وبگاه ها: شناخته شده ترین منابع وب

منابع وب آثاری هستند که توسط افراد مختلف در قالب های متنوع تولید و انتشار یافته اند. این منابع وبگاه ها، وبلاگ ها کتاب های الکترونیکی مقالات تمام نشریات الکترونیکی و... را در بر می گیرد. یکی از مهم ترین منابع تحت وب، وبگاه ها (1) هستند وبگاه مجموعه ای از صفحات وب است که دارای یک دامنه اینترنتی و به صورت مجموعه ای از صفحات مرتبط که داده هایی نظیر، متن صدا، تصویر و فیلم روی آن ها ارائه می شود، روی شبکه اینترنت قرار می گیرد صفحه وب به صورت سندی است که در قالب اچ تی ام ال (2) نوشته می شود و همواره با استفاده از پروتکل اچ تی تی پی (3) می توان به آن دسترسی پیدا کرد. مهم ترین قسمت یک وبگاه در واقع صفحه اصلی یا صفحه خانگی آن است وبگاه یک مؤسسه مرکز تجاری یا سازمان چهره او به سوی جهان و نقطه شروع بیش تر کاربران تلقی می شود نیلسن (2002) (4)

ارزیابی وبگاه ها

عواملی که مجموع قابلیت به رهبری از یک وبگاه را تعیین می کنند عبارتند از: محتوا، زبان، ساختار، طراحی جهت یابی عبور و قابلیت دسترسی روش های گوناگونی جهت ارزیابی قابلیت های بهره برداری از یک وبگاه وجود دارد که عبارتند از ارزیابی با مشارکت کاربر و ارزیابی بدون مشارکت کاربر روش ارزیابی در صورتی مفید خواهد بود که هم زبان و هم ساختار وبگاه مورد ارزیابی به سهولت توسط کاربران درک شود. (پل 2007) (5)

استانداردها و طرحهای ابر داده ای سازماندهی منابع وب

استانداردها و طرح های ابر دادهای مختلفی برای توصیف منابع اینترنتی تهیه شده اند که از انواع آن ها می توان به موارد ذیل اشاره نمود:

1. طرح ابر داده ای دابلین کور (6)

ص: 192

1- WebSites.

2- HTML

3- HTTP

4- J, Nielsen

5- Roswitha Poll

6- Dublin Core (DC)

2. مارک 21 (1)

3. آر دی اف (2)

4. طرح کدگذاری توصیف آرشیوی (3)

5. قالب ابر داده‌های خدمات مکان یاب اطلاعات دولتی (4)

6. قالب ابر داده ای طرح کدگذاری متن (5)

پردازش منابع الکترونیکی در سازمان اسناد و کتابخانه ملی ایران

در سازمان اسناد و کتابخانه ملی ایران برای پردازش و سازماندهی انواع منابع کتابی و غیر کتابی از ویرایش دوم قواعد فهرست نویسی انگلو امریکن (6)، استاندارد بین المللی توصیف کتاب شناختی (7) و استاندارد یونی مارک (8) استفاده می شود. کاربرد مارک (9) در سازماندهی منابع به طور جدی از سال 1385 آغاز شد. در این راستا از مارک ایران که شکل بومی سازی شده استاندارد یونی مارک است استفاده شده است. بر اساس این استاندارد منابع الکترونیکی در گروهی جداگانه و با کد 1 مشخص می شوند. در کتابخانه ملی ایران انواع کار برگه ها با توجه به این تقسیم بندی طراحی شده اند و هر کار برگه از 10 بلوک و فیلدهای اصلی و فرعی مربوط به توصیف هر منبع تشکیل شده است.

بر اساس طرح مقدماتی برای فهرست نویسی منابع الکترونیکی (2003) (10)، منابع الکترونیکی از لحاظ نوع دسترسی به دو به دو دسته منابع قابل «دسترسی مستقیم» (11) و منابع قابل «دسترسی از راه دور» (12) تقسیم می شوند (عبداللهی 1383). بنابراین وب گاه ها دسته ای از منابع الکترونیکی و از نوع «دسترسی از راه دور» به حساب می آیند منابع الکترونیکی در واحد سازماندهی منابع غیر کتابی اداره کل پردازش، سازماندهی می شوند. تاکنون بیش از 3600 منبع الکترونیکی از نوع «دسترسی مستقیم» در این واحد سازماندهی شده اند.

با توجه به مطالب پیش گفته به نظر می رسد که استانداردها و قواعد توصیف منابع الکترونیکی برای توصیف وبگاه ها نیز کاربردی باشد.

بیان مساله

وبگاه ها به عنوان نوعی از منابع اطلاعاتی وب در چرخه اطلاعات و دانش محسوب می شوند. با توجه به این که پردازش و سازماندهی منابع اطلاعاتی هسته مرکزی و فنی این چرخه به حساب می آید، لازم

ص: 193

1- Machin Readabel cataloging (MARC)

2- Encoded Archival description(EAD)

3- Resource Description Framework(RDF)

4- Government information locator sources (GILS)

Text Encoding Initiative(TEI)	-5
Anglo American Cataloging Rules, Second Edition (AACR2)	-6
International Standard Bibliographic Description (ISBD)	-7
Universal MARC	-8
Machine-Readable Cataloging	-9
Draft internet guidelines for cataloging electronic resources	-10
Direct access	-11
Remote access	-12

می نماید برای استفاده هر چه کامل تر از محتوای اطلاعاتی وبگاه ها سازمان اسناد و کتابخانه ملی ایران که مسئولیت استاندارد سازی و نظارت بر سازماندهی را بر عهده دارد از پردازش و سازماندهی وبگاه ها نیز غافل نشود با عنایت به اینکه کتابخانه ملی ایران در عرصه پردازش انواع منابع کتابخانه ای اعم از چاپی و الکترونیکی وارد شده است لازم است برای روزآمد سازی کنترل مدیریت و سازماندهی وبگاه ها نیز طرح و برنامه ای داشته باشد. این پژوهش بر آن است تا امکان پردازش وبگاه ها را به عنوان نوعی از منابع الکترونیکی مهم در عرصه های ملی و جهانی توسط سازمان اسناد و کتابخانه ملی ایران مورد پژوهش قرار دهد.

پیشینه پژوهش

پیشینه پژوهش در ایران:

حاجی زین العابدینی (1381) در پژوهشی به بررسی مشکلات اینترنت در زمینه سازماندهی و بازیابی اطلاعات پرداخته و دست نامه فهرست نویسی منابع اینترنتی را ارائه کرده است در تهیه دست نامه از قواعد فهرست نویسی انگلو امریکن قواعد به کار گرفته شده در طرح های فهرست نویسی منابع اینترنتی و الکترونیکی، اصول و قواعد فهرست نویسی منابع فارسی و نمونه های منابع اینترنتی فهرست نویسی شده، استفاده شده است. با استفاده از این دستنامه امکان فهرست نویسی منابع اینترنتی و سازماندهی اطلاعات دل خواه در اینترنت فراهم می شود هم چنین کاربرد های ایجاد شده است که بر اساس فیلدهای اطلاعاتی موجود در منابع، اینترنتی دستورالعمل های دست نامه نمونه های موجود در طرح هایی چون اینترکت، رهنمود های ایجاد مارک و فیلد های اطلاعاتی در راهنمای نرم افزار مارکایت طراحی شده است.

فتاحی حسن زاده (1385) به منظور مطالعه و ارزیابی شیوه های سازماندهی اطلاعات در وبگاه های کتابخانه های دانشگاهی پژوهشی انجام داده اند نتایج کلی پژوهش بیان گر آن است که در سازماندهی، اطلاعات صفحه اول وبگاه ها، شیوه دسته بندی بر اساس نوع خدمات بیش از سایر انواع رایج در حالی که برای سازماندهی سایر صفحات از شیوه های نسبتاً متنوعی استفاده می شود همچنین شیوه الفبایی عنوان و موضوع بیش ترین کاربرد و شیوه دایرکتوری (موضوعی سلسله مراتبی) کم ترین کاربرد را برای سازماندهی انواع منابع اطلاعاتی در وب سایت کتابخانه های دانشگاهی دارند.

پیشینه پژوهش در خارج از ایران:

کوچ (1) و همکارانش (1997) نقش طرح های رده بندی را در توصیف و بازیابی منابع اینترنتی مورد بررسی قرار داده اند. آن ها استفاده از این طرح ها را در سازماندهی محتویات وب گاه ها توصیه کرده اند، ولی به کاربرد شیوه های مختلف سازماندهی و نظریات کتابداران و کاربران در مورد چگونگی آن ها نپرداخته اند.

ویلیامسون (1997) (2) با تأکید بر ساختار دانش موجود در منابع اینترنتی در راستای سازماندهی دانش و

بازیابی اطلاعات، لزوم سازماندهی آن‌ها را بیان نموده است. در تحقیق وی، بررسی امکان کاربرد شیوه‌های متعارف سازماندهی برای سازماندهی محتوای اطلاعاتی وبگاه‌ها توصیه شده است.

وارد (2001) (1) در مقاله‌ای به تشریح اهمیت فهرست نویسی منابع اینترنتی پرداخته است. در ادامه فهرستی از فعالیت‌های انجام شده در کتابخانه‌های ایالات متحده در خصوص سازماندهی منابع اینترنتی را ارائه کرده است.

وایلر (2) و همکارانش (2008) در پژوهشی به ارزیابی هزینه‌های پردازش و سازماندهی منابع وب در کتابخانه و دانشگاه ملی کرواسی پرداخته‌اند نتایج این پژوهش نشان داده است که زمان پردازش منابع وبی با پردازش منابع چاپی یکسان است.

یانگه‌ی (3) (2011) در پژوهشی به بررسی فاصله بین ایجاد کنندگان منابع وب و ایجاد کنندگان فراداده‌های آن‌ها به منظور بهبود عناصر فراداده پرداخته است. یافته‌های پژوهش نشان‌گر این بوده است که رضایت کاربر تا حد بالایی به مفید بودن سهولت استفاده و دریافت اطلاعات و اثر بخشی عناصر اطلاعاتی بستگی دارد.

اهمیت پژوهش

سازماندهی و پردازش وبگاه‌ها در سازمان اسناد و کتابخانه ملی ایران می‌تواند از جنبه‌های زیر حائز اهمیت باشد:

1. ایجاد بانک اطلاعاتی وبگاه‌ها و غنی‌سازی کتاب‌شناسی ملی ایران؛
2. امکان جستجوی وبگاه‌ها با استفاده از قابلیت‌های جستجوی نرم‌افزار کتابخانه‌ای (رسا)؛
3. ایجاد نقاط دسترسی توصیفی و تحلیلی برای برآوردن نیازهای اطلاعاتی کاربران نهایی؛
4. ایجاد جامعیت و مانعیت در بازیابی پیشینه‌های کتاب‌شناختی وبگاه‌ها؛
5. استاندارد سازی و یکپارچه سازی پردازش توصیفی و تحلیلی وبگاه‌ها.

اهداف پژوهش

هدف کلی پژوهش دستیابی به فرایند پردازش و سازماندهی وبگاه‌ها با استفاده از قواعد و استانداردهای مورد استفاده در سازمان اسناد و کتابخانه ملی ایران است به این منظور اهداف فرعی زیر دنبال می‌شوند:

1. بررسی چگونگی دستیابی و دریافت اطلاعات وبگاه‌ها؛
2. توصیف کتاب‌شناختی وبگاه‌ها بر اساس استانداردها و قواعد مورد استفاده در سازمان اسناد و کتابخانه ملی ایران؛
3. دسترسی جامع و مانع به وبگاه‌ها با ایجاد نقاط بازیابی و تحلیل موضوعی مورد استفاده در سازمان اسناد و کتابخانه ملی ایران؛

Ward -1

Willer -2

Younghee -3

4. کاربرد نظام های موضوعی مورد استفاده در سازمان اسناد و کتابخانه ملی ایران برای نمایه سازی وبگاه ها

پرسش های اساسی

1. دستیابی و دریافت اطلاعات وبگاه ها از چه راه هایی امکان پذیر است؟

2. توصیف کتاب شناختی وبگاه ها بر اساس استانداردها و قواعد مورد استفاده در سازمان اسناد و کتابخانه ملی ایران چگونه است؟

3. دسترسی جامع و مانع به وبگاه ها با ایجاد نقاط بازیابی و تحلیل موضوعی مورد استفاده در سازمان اسناد و کتابخانه ملی ایران چگونه است؟

4. کاربرد نظام های موضوعی مورد استفاده در سازمان اسناد و کتابخانه ملی ایران برای نمایه سازی وبگاه ها چگونه است؟

تعاریف عملیاتی

امکان سنجی: منظور از امکان سنجی در این پژوهش بررسی طرح اولیه جهت ارائه و پیاده سازی فرایند پردازش و سازماندهی وبگاه ها در اداره کل پردازش و سازماندهی سازمان اسناد و کتابخانه ملی ایران است.

پردازش: منظور از پردازش در این پژوهش فرایند سازماندهی، منبع شامل دریافت، توصیف و تحلیل است.

جامعه آماری

جامعه آماری این پژوهش را 50 وبگاه تشکیل می دهند این جامعه مورد مطالعه در 20 رده موضوعی از وبگاه پارس ایندکس (1) و به صورت تصادفی انتخاب شده اند.

روش پژوهش و ابزار گردآوری داده ها

در این پژوهش دو روش کتابخانه ای و پیمایش توصیفی به کار گرفته شده است. جهت گردآوری اطلاعات از سیاهه واری و برای تجزیه و تحلیل داده ها از آمار توصیفی استفاده شده است. روایی سیاهه واری بر اساس تجربه پدید آوران مقاله در مدت بیش از یک دهه فعالیت تخصصی در حوزه سازماندهی منابع الکترونیکی و همچنین از طریق مشاوره با استادان به دست آمده است

یافته های پژوهش و پاسخ به پرسش های اساسی

پرسش اول: دست یابی و دریافت اطلاعات وبگاه ها از چه راه هایی امکان پذیر است؟

ص: 196

در این پژوهش جهت دستیابی و دریافت اطلاعات وبگاه‌ها دوروش مورد بررسی قرار می‌گیرد:

1. استفاده از فهرست وب گاه‌ها

برخی وب گاه‌ها و پایگاه‌های اطلاعاتی فهرست وب گاه‌های مختلف را بر اساس تقسیم بندی‌های موضوعی ارائه می‌دهند به این ترتیب یکی از راه‌های دست‌یابی به اطلاعات وب گاه‌ها، استفاده از فهرست وب گاه‌ها است. نمونه فهرست‌های وب گاه‌ها را می‌توان در وب گاه ایران (1) پارس ایندکس و الکسا (2) و غیره جستجو کرد.

در این پژوهش برای دسترسی به حجم نمونه از جامعه مورد مطالعه وبگاه پارس ایندکس انتخاب شد و امکان دست‌رسی به بسیاری از وب گاه‌ها امکان‌پذیر گردید.

2. ورود اطلاعات وبگاه‌ها در وبگاه سازمان اسناد و کتابخانه ملی ایران

منابع اینترنتی گروهی از منابع اطلاعاتی به حساب می‌آیند و کتابخانه ملی به لحاظ رسالت خود، لازم است مانند کتاب‌ها و سایر منابع اطلاعاتی اقداماتی در خصوص فراهم‌آوری منابع اینترنتی نیز به انجام رساند. در این راستا وب گاه سازمان اسناد و کتابخانه ملی ایران می‌تواند درگاهی جهت دست‌یابی به اطلاعات وب گاه‌ها باشد. بنابراین با طراحی «فرم الکترونیکی ورود اطلاعات وب گاه‌ها» و قرار دادن آن در وبگاه سازمان می‌توان به نشانی وب گاه‌ها و اطلاعاتی که برای پردازش آن‌ها مورد نیاز است، دست پیدا کرد.

جدول 1 فیلدهای پیشنهادی برای طراحی فرم الکترونیکی ورود وب گاه‌ها را نشان می‌دهد.

عکس

در این پژوهش جهت دستیابی و دریافت اطلاعات وبگاهها دو روش مورد بررسی قرار می‌گیرد:

۱. استفاده از فهرست وبگاهها

برخی وبگاهها و پایگاههای اطلاعاتی، فهرست وبگاههای مختلف را بر اساس تقسیم بندیهای موضوعی ارائه می‌دهند، به این ترتیب یکی از راههای دستیابی به اطلاعات وبگاهها، استفاده از فهرست وبگاهها است. نمونه فهرستهای وبگاهها را می‌توان در وبگاه ایران^۱، پارس ایندکس و الکسا^۲ و غیره جستجو کرد.

در این پژوهش برای دسترسی به حجم نمونه از جامعه مورد مطالعه، وبگاه پارس ایندکس انتخاب شد و امکان دسترسی به بسیاری از وبگاهها امکان پذیر گردید.

۲. ورود اطلاعات وبگاهها در وبگاه سازمان اسناد و کتابخانه ملی ایران

منابع اینترنتی گروهی از منابع اطلاعاتی به حساب می‌آیند و کتابخانه ملی به لحاظ رسالت خود، لازم است مانند کتابها و سایر منابع اطلاعاتی اقداماتی در خصوص فراهم‌آوری منابع اینترنتی نیز به انجام رساند. در این راستا وبگاه سازمان اسناد و کتابخانه ملی ایران می‌تواند در گامی جهت دستیابی به اطلاعات وبگاهها باشد. بنابراین با طراحی «فرم الکترونیکی ورود اطلاعات وبگاهها» و قرار دادن آن در وبگاه سازمان میتوان به نشانی وبگاهها و اطلاعاتی که برای پردازش آنها مورد نیاز است، دست پیدا کرد.

جدول ۱ فیلدهای پیشنهادی برای طراحی فرم الکترونیکی ورود وبگاهها را نشان می‌دهد.

جدول ۱. فیلدهای پیشنهادی برای طراحی فرم الکترونیکی ورود اطلاعات وبگاهها

نام فیلد	عنوان فارسی	عنوان به زبان دیگر	صاحب امتیاز	طراح رایانه ای	سال تولید	زمان آخرین ویرایش	زبان	حوزه موضوعی	کلید واژهها	نشانی اینترنتی	توضیحات
فیلد نوع	ضروری	عادی	ضروری	عادی	ضروری	ضروری	ضروری	ضروری	ضروری	ضروری	عادی

پرسش دوم: توصیف کتابشناختی وبگاهها بر اساس استانداردها و قواعد مورد استفاده در سازمان اسناد و کتابخانه ملی ایران چگونه است؟

جدول ۲ عناصر پیشنهادی توصیف وبگاهها را با استفاده از قواعد فهرست نویسی انگلومریکن در بخش منابع الکترونیکی و در بستر فرادادهای یونی مارک نشان می‌دهد.

1. <http://www.iran.ir/directory>

2. <http://www.alexa.com>

جدول 1 فیلدهای پیشنهادی برای طراحی فرم الکترونیکی ورود اطلاعات وبگاهها

پرسش دوم: توصیف کتاب شناختی وبگاهها بر اساس استانداردها و قواعد مورد استفاده در سازمان اسناد و کتابخانه ملی ایران چگونه است؟

جدول 2 عناصر پیشنهادی توصیف وبگاهها را با استفاده از قواعد فهرست نویسی انگلومریکن در بخش منابع الکترونیکی و در بستر فرادادهای یونی مارک نشان می‌دهد.

<http://www.iran.ir/directory> -1

<http://www.alex.com> -2

جدول ۲. فراوانی کاربرد قواعد انگلومریکن و استاندارد یونی مارک در توصیف کتابشناختی وبگاهها

درصد فراوانی	یونی مارک (بلوک)	یونی مارک (فیلد)	انگلومریکن (قواعد)	عناصر توصیف	ناحیه
%۷۶	اطلاعات توصیفی (۲)	Sar۰۰	B۹/۱	عنوان کامل	عنوان و پدیدآور
		Sbr۰۰	C۹/۱	وجه تسمیه عام: [منابع الکترونیکی]	
		Sdr۰۰	D۹/۱	عنوان به زبان دیگر	
		Ser۰۰	E۹/۱	دیگر اطلاعات عنوان	
		Sfr۰۰	F۹/۱	شرح پدیدآور	
%۶۸		Sar۰۵	B۹/۲	وضعیت ویراست	
بر اساس آخرین ویرایش قواعد انگلومریکن ناحیه ۳ (نوع انتشار) برای منابع الکترونیکی ذکر نمی‌شود					
%۷۰/۶	اطلاعات توصیفی (۲)	Sar۱۰	C۹/۴	محل تولید	چاپ و پدیدآور
		Ser۱۰	D۹/۴	ناشر و صاحب امتیاز	
		Sdr۱۰	F۹/۴	تاریخ تولید	
بر اساس قواعد انگلومریکن، ناحیه مشخصات ظاهری برای منابع الکترونیکی دسترسی از راه دور ذکر نمی‌شود					
%۵۱/۴	یادداشت‌ها (۳)	۳۰۰	۱B۹/۷	ماهیت و دامنه	یادداشت‌ها
		۳۰۰	۲B۹/۷	زبان	
		۳۰۴	۳B۹/۷	منبع عنوان کامل	
		۳۰۴	۶B۹/۷	شرح های پدیدآور	
		۳۰۵	۷B۹/۷	ویراست و تاریخچه کتابشناختی اثر	
		۳۰۶	۹B۹/۷	نشر، پخش و غیره	
		۳۰۷	۱۰B۹/۷	مشخصات ظاهری	
		۳۱۲	۴B۹/۷	عنوان های مرتبط	
		۳۱۴	۶B۹/۷	مسئولیت معنوی اثر	
		۳۲۷	۱۸B۹/۷	مندرجات	
		۳۳۰	۱۷B۹/۷	چکیده	
		۳۳۳	۱۴B۹/۷	مخاطبان	
		۳۳۶	۸B۹/۷	نوع منابع الکترونیکی	
۳۳۷	۸B۹/۷	روش دسترسی			
%۶۰	درصد فراوانی کل				

جدول 2 فراوانی کاربرد قواعد انگلومریکن و استاندارد یونی مارک در توصیف کتاب شناختی وب گاه ها

همان گونه که در جدول 2 مشاهده می شود بیش ترین کاربرد قواعد و استانداردهای مورد مطالعه برای توصیف وب گاه ها در ناحیه عنوان و پدیدآور با 76 درصد و کم ترین آن متعلق به ناحیه یادداشت با 51/4 درصد است. بنابراین در پاسخ به پرسش دوم می توان گفت امکان توصیف کتاب شناختی وب گاه ها اساس استانداردها و قواعد مورد استفاده در سازمان اسناد و کتابخانه ملی، ایران به میزان 60 درصد است.

پرسش سوم: دسترسی جامع و مانع به وب گاه ها با ایجاد نقاط بازیابی و تحلیل موضوعی مورد استفاده در سازمان اسناد و کتابخانه ملی ایران چگونه است؟

بخش دوم قواعد فهرست نویسی انگلوامریکن قواعد مربوط به نقاط بازیابی و تحلیل موضوعی را ارائه می دهد جدول 3 فیلد های بکار رفته در یونی مارک را بدین منظور نشان می دهد.

عکس

همان گونه که در جدول ۲ مشاهده می شود بیشترین کاربرد قواعد و استانداردهای مورد مطالعه برای توصیف وبگاهها در ناحیه عنوان و پدیدآور با ۷۶ درصد و کمترین آن متعلق به ناحیه یادداشت با ۵۱/۴ درصد است. بنابراین در پاسخ به پرسش دوم می توان گفت امکان توصیف کتاب شناختی وبگاهها بر اساس استانداردها و قواعد مورد استفاده در سازمان اسناد و کتابخانه ملی ایران، به میزان ۶۰ درصد است.

پرسش سوم: دسترسی جامع و مانع به وبگاهها با ایجاد نقاط بازیابی و تحلیل موضوعی مورد استفاده در سازمان اسناد و کتابخانه ملی ایران چگونه است؟
بخش دوم قواعد فهرست نویسی انگلوا مریکن قواعد مربوط به نقاط بازیابی و تحلیل موضوعی را ارائه می دهد. جدول ۳ فیلهای بکاررفته در یونی مارک را بدین منظور نشان می دهد.

جدول ۳. فراوانی کاربرد استاندارد یونی مارک در نقاط بازیابی و تحلیل موضوعی وبگاهها

درصد فراوانی	شماره بلوک	شماره فیلد	نام فیلد
%۶۱/۵	۵) آن های مرتبط	۵۱۰	عنوان اصلی به زبان دیگر
		۵۱۷	عنوان های گونه گون دیگر
		۵۳۲	عنوان گسترده
		۵۴۰	عنوان ترجمه شده
%۸۲	۶) موضوعی و تحلیلی	۶۰۰	نام شخص به منزله موضوع
		۶۰۱	نام تنالگان به منزله موضوع
		۶۰۵	عنوان به منزله موضوع
		۶۰۶	موضوع (اسم عام یا عبارت اسمی عام)
		۶۰۷	نام جغرافیایی به منزله موضوع
%۵۳/۳	۷) مقوری اثر	۷۰۲	نام شخص به منزله شناسه افزوده
		۷۱۰	نام تنالگان به منزله سرشناسه
		۷۱۲	نام تنالگان به منزله شناسه افزوده
%۱۰۰	۸) کاربردی	۸۵۶	نشانی اینترنتی
%۷۰/۴	درصد فراوانی کل		

با توجه به جدول ۳، بیشترین کاربرد استاندارد یونی مارک مربوط به بلوک ۸ با ۱۰۰ درصد است، بعد از آن بلوک ۶ با ۸۲ درصد و بلوک ۵ با ۶۱/۵ درصد قرار دارد. کمترین کاربرد نیز در بلوک ۷ با ۵۳/۳ درصد مشاهده می شود. بنابراین در پاسخ به پرسش سوم میتوان گفت امکان دسترسی جامع و مانع

جدول 3 فراوانی کاربرد استاندارد یونی مارک در نقاط بازیابی و تحلیل موضوعی وب گاه ها

با توجه به جدول، بیش ترین کاربرد استاندارد یونی مارک مربوط به بلوک 8 با 100 درصد است بعد از آن بلوک 6 با 82 درصد و بلوک 5 با 61/5 درصد قرار دارد کم ترین کاربرد نیز در بلوک 7 با 53/3 درصد مشاهده می شود بنابراین در پاسخ به پرسش سوم می توان گفت امکان دسترسی جامع و مانع

به وب گاه ها با ایجاد نقاط بازیابی و تحلیل موضوعی مورد استفاده در سازمان اسناد و کتابخانه ملی ایران حدود 70 درصد است.

پرسش چهارم: کاربرد نظام های موضوعی مورد استفاده در سازمان اسناد و کتابخانه ملی ایران برای نمایه سازی وب گاه ها چگونه است؟

در کتابخانه ملی، ایران منابع غیر کتابی با استفاده از زبان کنترل شده نمایه سازی می شوند. جدول 4 نظام های موضوعی مورد استفاده را نشان می دهد لازم به ذکر است علاوه بر نظام های موضوعی زیر، از سر عنوان های موضوعی کتابخانه کنگره (1)، سایر اصطلاحنامه ها و واژه نامه های موضوعی، دایره المعارف ها و بانک های اطلاعاتی عمومی و موضوعی جهت مستندسازی توصیف گر ها و پیشنهاد آن ها به نظام های موضوعی پیش گفته استفاده می شود

عکس

به وبگاه‌ها با ایجاد نقاط بازبایی و تحلیل موضوعی مورد استفاده در سازمان اسناد و کتابخانه ملی ایران حدود ۷۰ درصد است.

پرسش چهارم: کاربرد نظام‌های موضوعی مورد استفاده در سازمان اسناد و کتابخانه ملی ایران برای نمایه‌سازی وبگاه‌ها چگونه است؟

در کتابخانه ملی ایران، منابع غیرکتابی با استفاده از زبان کنترل‌شده نمایه‌سازی می‌شوند. جدول ۴ نظام‌های موضوعی مورد استفاده را نشان می‌دهد. لازم به ذکر است علاوه بر نظام‌های موضوعی زیر، از سرعنوان‌های موضوعی کتابخانه کنگره^۱، سایر اصطلاحنامه‌ها و واژه‌نامه‌های موضوعی، دایره‌المعارف‌ها و بانک‌های اطلاعاتی عمومی و موضوعی جهت مستندسازی توصیفگرها و پیشنهاد آنها به نظام‌های موضوعی پیشگفته استفاده می‌شود.

جدول ۴. فراوانی کاربرد نظام‌های موضوعی در نمایه‌سازی وبگاه‌ها

نظام موضوعی	حوزه موضوعی	درصد فراوانی
اصطلاحنامه فرهنگی فارسی (اصفا)	علوم انسانی	٪۶۶
اصطلاحنامه پزشکی فارسی	پزشکی	٪۱۰۰
اصطلاحنامه های علوم	فنی - مهندسی	٪۱۰۰
سرعنوان های موضوعی فارسی	همه علوم	٪۷۲
درصد فراوانی کل		٪۸۴/۵

بر اساس جدول ۴، هر کدام از نظام‌ها در حوزه موضوعی خود در جامعه مورد مطالعه، بررسی شده‌اند. بنابراین در پاسخ به پرسش چهارم می‌توان نتیجه گرفت که کاربرد نظام‌های موضوعی مورد استفاده در سازمان اسناد و کتابخانه ملی ایران برای نمایه‌سازی وبگاه‌ها بالای ۸۴ درصد امکان‌پذیر است.

بحث و نتیجه‌گیری

یافته‌های پژوهش نشان می‌دهند پردازش وبگاه‌ها در سازمان اسناد و کتابخانه ملی ایران، بیش از ۶۵ درصد امکان‌پذیر است. بر اساس نتایج این پژوهش برای طراحی فرایند پردازش وبگاه‌ها می‌توان فرم الکترونیکی ثبت وبگاه‌ها را طراحی کرد و آن را در وبگاه سازمان دسترسی‌پذیر نمود. در برخی موارد نیز با استفاده از فهرست وبگاه‌ها می‌توان نوعی فراهم‌آوری را انجام داد. در مرحله بعدی با استفاده از قواعد فهرست‌نویسی انگلوا امریکن در بستر فراداده‌های یونی‌مارک و بومی‌سازی آن و با استفاده از کاربرگه

1. (LCSH) Library of Congress Subject Headings

جدول 4 فراوانی کاربرد نظام‌های موضوعی در نمایه‌سازی وب‌گاه‌ها

بر اساس جدول 4، هر کدام از نظام‌ها در حوزه موضوعی خود در جامعه مورد مطالعه بررسی شده‌اند. بنابراین در پاسخ به پرسش چهارم می‌توان نتیجه گرفت که کاربرد نظام‌های موضوعی مورد استفاده در سازمان اسناد و کتابخانه ملی ایران برای نمایه‌سازی وب‌گاه‌ها بالای 84 درصد امکان‌پذیر است.

یافته های پژوهش نشان می دهند پردازش وب گاه ها در سازمان اسناد و کتابخانه ملی ایران، بیش از 65 درصد امکان پذیر است. بر اساس نتایج این پژوهش برای طراحی فرایند پردازش وب گاه ها می توان فرم الکترونیکی ثبت وب گاه ها را طراحی کرد و آن را در وب گاه سازمان دسترس پذیر نمود. در برخی موارد نیز با استفاده از فهرست وب گاه ها می توان نوعی فراهم آوری را انجام داد در مرحله بعدی با استفاده از قواعد فهرست نویسی انگلوا امریکن در بستر فرادادهای یونی مارک و بومی سازی آن و با استفاده از کاربرگه

ص: 200

منابع الکترونیکی اقدام به توصیف و تحلیل وبگاه‌ها نمود سپس آن‌ها را بر اساس نقاط بازیابی مختلف دسترس پذیر ساخت. با توجه به نتایج پژوهش، نمایه‌سازی و تحلیل موضوعی وبگاه‌ها نیز با استفاده از نظام‌های موضوعی پیش‌گفته امکان پذیر خواهد بود.

در مقایسه نتایج این پژوهش با پژوهش‌های پیشین می‌توان خاطر نشان کرد نمونه‌سازماندهی وبگاه‌ها در کتابخانه‌ها در پژوهش «وایلر» و همکارانش و همچنین «وارد» مورد بررسی قرار گرفته است. پردازش وبگاه‌ها از لحاظ هزینه و رضایت کاربران همان‌گونه که «وارد»، «وایلر» و همکارانش و هم چنین «یانگهی» در پژوهش خود اشاره کرده‌اند از اهمیت و ارزش والایی برخوردار است. فتاحی و «حسن زاده» و همچنین «کوچ» و همکارانش نیز بر اعمال رده‌بندی و فهرست موضوعی، جهت دسترسی به وبگاه‌ها تأکید کرده‌اند. همچنین استفاده از قواعد و استانداردهای مرسوم سازماندهی از جمله قواعد فهرست نویسی انگلوامریکن و مارک برای سازماندهی وبگاه‌ها در پژوهش «حاجی زین العابدینی»، «ویلیامسون» و «کوچ» و همکارانش مورد بررسی و تأیید قرار گرفته است.

به طور کلی از این پژوهش می‌توان نتیجه گرفت که استاندارد سازی موضوع سازماندهی و پردازش وبگاه‌ها به عنوان نوعی از منابع الکترونیکی و بومی سازی آن در سازمان اسناد و کتابخانه ملی ایران امکان پذیر است.

پیشنهاد‌های برخاسته از پژوهش

پس از انجام این پژوهش و با عنایت به یافته‌های آن پیشنهاد می‌شود سازمان اسناد و کتابخانه ملی ایران نسبت به انجام موارد زیر که جنبه عملیاتی دارند اقدام نماید.

1. طراحی و ایجاد سامانه ورود اطلاعات وبگاه‌ها در وبگاه سازمان اسناد و کتابخانه ملی ایران جهت فراهم‌آوری و پردازش آن‌ها؛
2. ارائه آموزش‌های تخصصی در حوزه پردازش و سازماندهی وبگاه‌ها به نیروی انسانی و کارشناسان ذیربط؛
3. برنامه ریزی در جهت اشاعه اطلاعات وبگاه‌ها به صورت جامع و مانع.

منابع

پارس ایندکس <http://parsindex.com/default.aspx> (دسترسی در 1391/7/20)

پژوهشگاه علوم و فناوری اطلاعات ایران 1390 اصطلاح‌نامه‌های علوم <http://thesauri.irandoc.ac.ir> (دسترسی در 1391/7/20)

پل روزیتا. 1390. ارزیابی وبگاه کتابخانه‌ای پژوهش‌های آماری و معیارهای کیفیت ترجمه رضا خانی پورگزیده مقالات ایفلا 2007 دوربان آفریقای جنوبی، 23-19 اوت 2007 تهران نشر کتابدار

حاجی زین العابدینی، محسن . 1381. بررسی مسائل فهرست نویسی منابع اینترنتی و ارائه دست نامه پیشنهادی برای کتابخانه های ایران پایان نامه کارشناسی ارشد دانشگاه علوم پزشکی و خدمات بهداشتی درمانی ایران دانشکده مدیریت و اطلاع رسانی پزشکی.

سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران بانک مستندات <http://opac.nlai.ir> (دسترسی در 1391/7/20).

عبداللہی، مریم 1383 ارزیابی وضعیت فهرست نویسی منابع الکترونیکی بر اساس قواعد انگلوماریکن در کتابخانه های دانشگاهی و آرشیوهای موجود در شهر تهران پایان نامه کارشناسی ارشد. دانشگاه آزاد، اسلامی واحد تهران شمال

علی محمدی، داریوش، 1381. ابر پایگاه های اطلاعاتی در شبکه جهانی. وب فصلنامه کتاب. 49(1): 93-100

فتاحی، رحمت الله. حسن زاده محمد. 1385 نظر سنجی از کتابداران متخصص پیرامون شیوه های سازماندهی اطلاعات در وب سایت کتابخانه های دانشگاهی گزارشی از مرحله دوم یک طرح پژوهشی فصلنامه کتابداری و اطلاع رسانی 36 (4): 6-30.

کمیته ملی مارک ایران 1381 مارک ایران تهران: کتابخانه ملی جمهوری اسلامی ایران

نشاط، نرگس. 1382 چالش های سازماندهی موضوعی منابع وب اطلاع شناسی. 1: 37-54.

IFLA. 2008. UNIMARC concise bibliographic format. Available at: <http://archive.ifla.org/VI/8/unimarc-concise-bibliographic-format-2008.pdf> [September. 2012]

Koch, Traugott et al. 1997. The role of classification schemes in Internet resource description and discovery. [Available at: www.ub.lu.se/desire/radar/reports/D3,2,3/class-v10.html] [September. 2012]

Nielsen, J. 2002. Top 10 guidelines for homepage usability. Available at: <http://www.useit.com/alertbox/20020512.htm> [September. 2012]

Steering Committee of American Library Association, et al.. 2002. Anglo American Cataloging Rules. Joint Chicago: American Library Association

Ward, Diane. 2001. Internet resource cataloging: the SUNY Buffalo Libraries' response. OCLC Systems Services, 17 (1): 19-26

Willer, Mirna. Buzina, Tanja. Holub, Karolina. Zajec, Jasenka. Milinovic, Miroslav. Topolscak, Nebojsa. 2008. Selective Archiving of Web Resources: A Study of Processing Costs. Program: Electronic Library and Information Systems. v42. n4: p341-364

Williamson N. J.1997. Knowledge structures and the Internet. Knowledge organization for information retrieval: proceedings of the sixth international study conference on classification research. University College London 16-18 June: pp23-27

Younghee Noh. 2011. A study on metadata elements for web-based reference resources system developed through usability testing. Library Hi Tech, 29(2):242-265

ص: 203

محاسبه هزینه‌ها دانسته‌های کنونی ما را در خصوص ایجاد آرشیو تحت وب و اینکه این نوع آرشیو در مقایسه با شیوه فراهم‌آوری منابع چاپی، فعالیتی پرهزینه بوده را تأیید می‌کند. مقاله حاضر اطلاعات سودمندی را در خصوص ایجاد یک آرشیو گزینشی منابع تحت وب و چگونگی ارتقای برنامه‌های نرم‌افزاری ارزیابی مجموعه‌های فراهم‌آوری شده را در اختیار خواننده قرار می‌دهد. و در مقاله راهکارهایی ارائه می‌شود که منجر به افزایش دانش در زمینه بهره‌گیری از فناوری‌های مربوط به ایجاد آرشیو تحت وب، برای افزایش کارآمدی تمامی پروژه‌های مربوط به ایجاد آرشیو تحت وب می‌شود. یکی از این پروژه‌ها که در مقاله بررسی شده است، آرشیو کتابخانه ملی استرالیاست که اخیراً با استفاده از رویکرد گزینشی، موفق شده که هزینه‌های سرانه خود را در یک بازه زمانی کاهش دهد.

طی دهه، گذشته تعدادی اندکی - اما - اما رو به رشد - از کتابخانه های ملی اقدام به تدوین برنامه هایی در زمینه ایجاد آرشیو تحت وب کرده اند. برنامه هایی که یک یا بیش از چهار رویکرد اصلی زیر را در برداشته اند:

1. ایجاد آرشیو گزینشی، برای مثال، آرشیوهای کتابخانه ملی، کانادا، ژاپن و استرالیا؛

2. تهیه و تدارک دامنه وب در کل کشور به طور دوره ای، برای مثال نمونه های انجام گرفته توسط آرشیو های کشور های اسکاندیناوی، از جمله سوئد؛

3. گردآوری موضوعی، برگرفته از کتابخانه کنگره مجموعه مینروا از منتخب های سال 2000، 2002، و 11 سپتامبر 2001؛ و

ص: 205

مجموعه های واسپاری نظیر STORS کتابخانه ایالتی تاسمانی و واسپاری الکترونیکی (1) کتابخانه ملی هلند

کتابخانه ملی، استرالیا به دلیل مزایای موجود در رویکرد گزینشی این روش را انتخاب کرد این دلایل عبارت اند از:

• هر یک از اقلام موجود در آرشیو کتابخانه از لحاظ کیفیت مورد ارزیابی قرار گرفته و تا سر حد امکان با بهره گیری از قابلیت های فنی موجود در این زمینه کاربرد پذیر می شوند.

• هر یک از اقلام موجود در آرشیو را می توان به طور کامل فهرست نویسی، و به بخشی از کتاب شناسی ملی تبدیل کرد و داده های کتاب شناختی را به اشتراک گذاشت هم چنین در فهرست کتابخانه، اطلاعات کتاب شناختی منابع تحت وب با منابع دیگر ادغام شده و کاربران می توانند به صورت یکجا به تمامی منابع دسترسی پیدا کنند.

• هر کدام از اقلام موجود در وب می تواند به سرعت از طریق محیط وب برای مخاطبان خود دسترس پذیر شود؛ زیرا پیش تر تلاش شده اجازه انجام این کار از ناشران اخذ شود.

• ویژگی های هر یک از منابع موجود در آرشیو و شیوه طبقه بندی آن ها برای مدیران مجموعه شناخته شده است. این امر در نخستین گام توانایی ما را در توسعه روش ها و ابزار های گردآوری منابع ذخیره سازی و دسترس پذیری به آن ها افزایش می دهد افزون بر این شناخت و آگاهی، مدیران مجموعه را بر آن می دارد که سیاست ها و راهبرد های حفاظت و نگهداری منابع را به گونه ای اتخاذ کنند که امکان دسترس پذیر سازی منابع برای مدت زمان طولانی فراهم شود.

• سایت هایی را که برای روبات های جست و جوگر موجود دسترس پذیر نیستند نیز می توان شناسایی و گردآوری نمود و با استفاده از روش های دیگر- برای مثال بر حسب نام ناشر - مرتب کرد. این امر، نام های تجاری - که نیازمند رمز عبور ناشر هستند- و پایگاه داده ها را نیز در بر می گیرد.

با وجود مزایای فوق هر یک از رویکرد های موجود در زمینه ایجاد آرشیو تحت وب، معایبی نیز دارد. رویکرد گزینشی نیز از این قاعده مستثنا نیست این، امر به کارمندان کتابخانه بستگی دارد که در محیطی مملو از اطلاعات کاملاً جدید به ایفای وظیفه می پردازند، عملکرد آن ها در این زمینه، باید به گونه ای باشد که در نهایت به قضاوت در خصوص آن چه که در آینده باید مورد تحقیق و تفحص قرار گیرد، بی انجامد و عرصه تحقیقات آینده را روشن سازد. رویکرد گزینشی همچنین به استخراج منابع، جدا از محتوا و پیوند میان آن ها با دیگر منابع خارجی می پردازد. از نظر توسعه و مدیریت آرشیو، بزرگ ترین نقطه ضعف رویکرد آرشیوی این است که روند کار فشرده و پر زحمت است و هزینه انجام کار به ازای هر منبع که قرار است آرشیو شود، بالاست.

ایجاد آرشیو گزینشی منابع تحت وب

اشاره

ایجاد آرشیو گزینشی منابع تحت وب (2)

بررسی هزینه های مربوط به فراهم آوری منابع تحت و در کتابخانه ملی استرالیا

تاریخچه

کتابخانه ملی استرالیا در 1996، به صورت گزینشی و آزمایشی اقدام به ایجاد آرشیوهای تحت وب کرد.

ص: 206

e-Depot -1

Selective Archiving of Web Resources: A Study of Acquisition Costs at the National Library of Australia -2

Margaret Philips, Director, Digital Archiving, National Library of Australia -3

4- استادیار کتابداری و اطلاع رسانی و دانشیار گروه کتابداری و اطلاع رسانی دانشگاه آزاد اسلامی واحد علوم و تحقیقات تهران

اقدامی که در تاریخ خود حرکتی بسیار نو محسوب می شد و کمتر نوشته ای بدان پرداخته بود. حتی کسی نبود که بتوان از او آموخت. ما در مورد ایجاد آرشیوگان تحت وب اطلاعات کافی نداشتیم تا بتوانیم در این خصوص اقدام به تدوین طرح یا راهبردهای مناسبی جهت برآورد هزینه های مربوط به این کار کنیم. تنها در سایه برداشتن گام های کوچک عملی و با عنوان خودآموزی توانستیم به راهمان ادامه داده و رشد کنیم هیچ گونه کمک مالی خاصی برای انجام این کار جدید دریافت نمی کردیم، فقط مجبور بودیم آن دسته از کارکنان مجموعه سازی را که جدا از وظایف کاملاً سنتی، خود، به کار جدید ایجاد آرشیوگان تحت وب علاقه مند بوده و از خود استعدادهایی در این زمینه بروز می دادند، برای انجام این کار جدید تربیت کنیم برای مثال شخصی در بخش فناوری اطلاعات باید بخشی از وقت خود را صرف این می کرد که چگونه می تواند منابع را از وب گاه های ناشران دریافت و اطلاعات آن ها را بر روی سرور کتابخانه نگهداری کند ما برای در اختیار گرفتن و مدیریت فایل های دریافت شده، تنها مجبور به استفاده از نرم افزارهایی بودیم که به طور رایگان موجود بود با اعمال چنین تدابیری نه تنها هزینه های اولیه ایجاد آرشیو تحت وب کاملاً اندک، بلکه عمدتاً در بودجه مربوط به کارکنان پنهان می شد.

با گذشت 9 سال عملیات مربوط به مراتب پیچیده تر شدند کتابخانه استرالیا نرم افزار پاندورا (1) یعنی آرشیو تحت وب استرالیا (2) را راه اندازی کرد که به مثابه آرشیوی کارآمد، وارد طرح شراکت با 9 کتابخانه دیگر استرالیا و نیز دیگر سازمان های مسئول گردآوری منابع فرهنگی، یکی پس از دیگری شد. در مجموع می توان گفت که این امر روندی رو به رشد و توسعه داشت، اما نهادها را با افزایش هزینه ها مواجه ساخت حجم افزونتر فعالیت های مربوط به ایجاد آرشیو و نیاز آن به پشتیبانی از مشارکت کنندگانی که از مکان های دور در این طرح شرکت می کردند، ضرورت ایجاد زیر ساخت های پیچیده فنی متشکل از سیستم ارسال و تحویل مدارک، مدیریت آرشیو، و ذخیره سازی را ایجاب می کرد این، امر به مفهوم وجوب هزینه های ضروری در امر توسعه بود، هر چند که این هزینه ها همچنان از طریق بودجه های موجود پرسنلی برای توسعه مجموعه سازی و بخش فناوری اطلاعات فراهم می شد.

کتابخانه مذکور نمی توانست نظام مدیریت آرشیو مناسبی برای خریداری پیدا کند، بنابراین، نظام ایجاد آرشیو دیجیتالی پاندورا (پانداس) (3) به همراه نظام تحویل مدارک پاندورا، پا به عرصه وجود

گذاشتند. همچنین کتابخانه اقدام به خرید سیستم ذخیره سازی شیء دیجیتالی (دی او اس اس) کرد که زمینه اشتراک منابع را میان پاندورا و دیگر مجموعه های دیجیتالی کتابخانه را فراهم می کرد.

یکی از چیزهایی که در 9 سال گذشته تغییر نکرده این است که: کتابخانه هنوز هیچ بودجه اضافی برای انجام این فعالیت ها در اختیار ندارد به طرز عجیب شواهد موجود بر سودآور بودن چنین اقدامی دلالت دارد؛ لذا چنین مسئله ای ما را بر آن می دارد که بودجه ای را برای تحقق این امر از محل بودجه دولت استرالیا که هر ساله رقم آن رو به تزاید است تعریف نماییم با این حال هیچ گونه طرح کوتاه

ص: 207

مدت خاصی که بتواند بودجه ریزی این فعالیت را به گونه ای توجیه کند تا هنگام [طرح] اتمام به ثبات این فعالیت بی انجامد و آن را در حالت بی ثباتی و ناپایداری رها نکند، وجود ندارد.

آرشیو پاندورا در حال حاضر

آرشیو پاندورا مجموعه قابل توجهی از نشریات برخط و وب گاه های استرالیاست که توسط کتابخانه ملی این کشور و شرکای آن ایجاد شده و گسترش یافته. است آرشیوی که در کتابخانه ملی کانبرا (1)، به شیوه ای متمرکز ذخیره، اداره، و حفظ و نگهداری می شود. از 30 آوریل 2005، که آرشیو مشتمل بر 8235 عنوان، بود، تا کنون نرخ متوسط رشد مجموعه سالیانه در حدود 2400 عنوان می باشد. این عناوین شامل فایل مجزایی از اسناد متنی نظیر فایل های پی دی اف (2)، یا اقلام تحت، وب، نظیر وب گاه های حجیم و پیچیده، متشکل از هزاران فایل در اشکال و فرمت های، مختلف شامل، متن، صدا، تصویر و یا ویدئوست که بر پایه نظمی مشخص و به منظور ایجاد «شاخصی» جدید از عناوین آرشیوی ای که در حال حاضر رقم آن ها در آرشیو موجود به معادل 16/736 مورد بالغ می گردد گردآوری شده است. هدف این بود که عناوین منتخب به شیوه ای ایمن و صحیح در آرشیو ذخیره سازی و به طور مستمر برای همگان دسترس پذیر شود.

آرشیو مزبور شامل هر دو نوع انتشارات برخط ایستا و پویا (3) بوده و وب گاه ها را نیز در بر می گیرد، و نشان دهنده طیف گسترده ای از انواع انتشارات و فرمت های استفاده شده توسط ناشران و پدید آورندگان بر روی وب است. افزون بر این نشریات برخط و وب گاه هایی را که هم اکنون در محیط زنده وب ناپدید گردیده و در هیچ کجای دیگر نیز در دسترس نمی باشند نیز در بر می گیرد.

بسیاری از عناوین موجود در این آرشیو می تواند با اتصال به اینترنت به راحتی و آزادانه در دسترس تمامی افراد در هر نقطه از جهان قرار گیرد تقریباً پنج سال است - اندکی کمتر - که به دلایل تجاری دسترسی به بخش بسیار کوچکی از آرشیو محدود شده است. با وجود این، اطلاعات مربوط به این عناوین دارای این محدودیت را می توان بر روی کامپیوتری واحد در اتاق مطالعه اصلی کتابخانه مشاهده کرد.

دسترسی به محتویات آرشیو را می توان یا از طریق ایجاد پیوندی مناسب به پیشینه فهرست نویسی منبعی خاص یا از طریق موضوع و فهرست عناوین موجود در وب گاه پاندورا امکان پذیر کرد. موتور های جست و جوی متعارف مانند گوگل و یاهو، به نمایه سازی منابع آرشیوی بر حسب آن چه که در عنوان آن ها آمده است می پردازند.

محدوده وظایف

وظایفی که کارکنان کتابخانه ملی و دیگر شرکای آن ها به عنوان بخشی از برنامه ایجاد آرشیو تحت وب پاندورا انجام می دهند متأثر از مجموعه عوامل متعددی نظیر تصمیمات مربوط به سیاست گذاری ها و شرایط محیطی است؛ مجموعه سیاست ها و عواملی که نه به شکلی، حداقل بلکه به شیوه ای کاملاً مؤثر بر

ص: 208

افراد شرکت کننده در طرح پاندورا به بالاترین شکل بر حفظ اصل به طرفه العینی به نتایج مناسب (رسیدن خصوصیات ظاهری و کارآمد پذیری منبع) در باب منبع منتشره و یا وب گاه، و نیز توجه همزمان بر محتویات آن تأکید می ورزند به محض آن که فرد گردآوری کننده نسخه ای از یک منبع را بر روی سرور کتابخانه ملی قرار دهد افراد همکار در سازمان های دیگر نیز نسخه مزبور را جهت کسب اطمینان از جامعیت و کارآمدی آن و قبل از ارسال آن منبع به آرشیو جهت استفاده عموم، مورد بازبینی قرار می دهند. این تضمین کیفیت فرآیندی بسیار وقت گیر و در نتیجه گران قیمت و پرهزینه است که در واقع، گران ترین جنبه فرآیند فراهم آوری منابع نیز محسوب می شود.

هر عنوان در آرشیو بر مبنای رکوردهای موجود در کتابخانه ملی فهرستگان پیوسته کتابخانه های، دیگر و بر مبنای استفاده از اطلاعات پایگاه کتابشناختی ملی (فهرستگان مشترک اطلاعات کتاب شناختی منابع مربوط به بیش از 850 کتابخانه استرالیا به همراه فراهم آوری امکان دست یابی به منابع توسط خدمات مرورگر کین تکا (1)، فهرست نویسی گردیده است) تصمیم در اتخاذ چنین سیاستی از آن جهت صورت گرفته که امکان شناخت و استخراج منابع برخط یکپارچه سازی شده با دیگر منابع کتابخانه ای، امری بسیار مهم تشخیص داده شده است با وجود، این محرز است که انجام چنین اقداماتی بر هزینه های مربوط به ایجاد آرشیو منابع تحت وب می افزاید.

عامل مهم و یاری کننده دیگری که به صرفه جویی زمانی و استفاده بهینه از اوقات کاری کارکنان منجر می شود شرایط مطرح در قوانین و اسپاری کشور استرالیا است. در بسیاری از ایالت های استرالیا، حتی آن هایی که مشترک المنافع هستند نیز قانون و اسپاری منابع و از آن جمله منابع برخط همچنان مهجور مانده و به تصویب نرسیده است تنها در بخش های شمالی این کشور اخیراً قانونی به تصویب رسیده است که به وضوح امکان واسپاری منابع برخط را میسر می سازد بدان معنا که تمامی کتابخانه های همکار در این طرح از جمله کتابخانه ملی باید قبل از تهیه نسخه ای از یک منبع نسبت به کسب اجازه از ناشر آن، اقدام سپس آن را برای استفاده همگان در آرشیو دسترس پذیر کنند.

نیروی انسانی و افراد شاغل در پاندورا

ایجاد آرشیو پاندورا اجرا و گسترش برنامه های کاربردی وابسته به آرشیو حفظ و نگهداری سیستم های مربوط، تخصص گرایی و برنامه ریزی به منظور محافظت بلند مدت از آن مستلزم به کارگیری نیروهای کاری شش واحد در دو بخش کتابخانه ملی است

کارکنان شاغل در آرشیو دیجیتال بخش مجموعه سازی مسئول انتخاب منابع موجود در آرشیو و محتوای آن می باشند واحد برنامه های کاربردی بخش فناوری اطلاعات نیز مسئول توسعه زیر ساخت های فنی آرشیو هستند بخش مربوط به خدمات وبگاه مسئولیت توسعه رابط کاربر را برای هر دو سیستم پاندورا و پانداس بر عهده دارد و نظام های حمایتی سیستم های تجاری خود عاملی جهت حفظ تولیدات

آزمایش آموزش و ارزیابی محسوب می شوند و موجب می گردند اهداف مذکور به خوبی اجرا و محقق شود. افزودن بر، این کارکنان واحد خدمات حفاظتی بخش مجموعه سازی مسئولیت دسترس پذیر سازی طولانی مدت محتوای آرشیو را بر عهده دارند.

در سال 2004، در زمان تخمین و برآورد هزینه های زیر سهم کمک های کتابخانه ملی به آرشیو پاندورا در قالب استفاده از نیروهای تمام وقت به شرح سطوح زیر بود:

شعبه آرشیو دیجیتال:

• 1 نفر نیروی تمام وقت در سطح مدیر - کتابدار

• 2 نفر نیروی تمام وقت در سطح یک نفر سرپرست یک نفر کارشناس پروژه های خاص -

• کتابداران

• 4 نفر نیروی تمام وقت در سطح کارکنان ستادی - کتابداران؛ و

• 17 نفر نیروی تمام وقت در سطح کارشناس رفع مشکلات فنی - دارای 5 سال سابقه کار در حوزه فناوری اطلاعات

• 2 نفر مدیر تمام وقت که 10 درصد از وقت خود را هر کدام صرف توسعه و نگهداری سیستم کنند؛

و

• 2 نفر مدیر تمام وقت که 25 درصد از وقت خود را هر کدام صرف توسعه و نگهداری سیستم نمایند. همان طور که در بالا ذکر شد نیروی انسانی شاغل در بخش حفاظت از منابع آرشیوی، به صورت تمام وقت در این طرح حضور داشتند، با وجود این هزینه های مربوط در این برآورد منظور نگردیده است.

هزینه دستیابی به وب گاه ها و نشریات برخط

گر چه برای کتابخانه کاملاً مشخص بود که هزینه سرانه ایجاد واحد آرشیو منابع برخط در مقایسه با منابع چاپی نظیر کتاب ها و پیاپی ها بالا می باشد؛ با این حال، تا همین اواخر هیچ گونه اطلاعات جامع و دقیقی در خصوص هزینه های ایجاد آرشیو تحت وب در اختیار نداشتیم. در سال 2004، کتابخانه تصمیم گرفت تا هزینه های مربوط به فعالیت های واسپاری منابع پیاپی و تک نگاشت ها را برآورد کند. افزودن بر این، از آن جا که به نظر می رسید گردآوری منابع برخط استرالیا نیز امتداد بسط همان مسئولیت های مذکور باشد، تصمیم گرفته شد هزینه اجرای این کار نیز به همان روال برآورد شود.

محدوده برآورد هزینه

مطالعات مربوط به برآورد هزینه به بررسی و مطالعه هزینه های مربوط به فراهم آوری منابع کتابخانه ملی به عنوان یک «نمونه» (1) پرداخته سپس آن را به هزینه های پاندورا می افزاید افزودن بر این مرز میان هزینه هایی که باید لحاظ یا حذف شوند نیز مشخص می گردد.

هزینه های مستقیم عبارت بودند از:

• هزینه های کارکنان بخش آرشیو دیجیتال کتابخانه ملی؛

• هزینه های اداری بخش آرشیو دیجیتال از قبیل، مسافرت، آموزش شرکت در کنفرانس و

• تجهیزات اداری (شامل هزینه های فرد تأمین کننده)؛ و

• هزینه های مربوط به توسعه زیر ساخت ها و نگهداری از آن ها نظیر هزینه کارکنان بخش فناوری اطلاعات و خرید تجهیزات سخت افزاری و نرم افزاری

هزینه هایی که باید حذف شوند عبارت بودند از:

• هزینه های غیر مستقیم مانند تأمین ایستگاه های کاری کارکنان

• نورپردازی و محافظت از ساختمان و

• هزینه های مربوط به محافظت از محتوای منابع موجود در آرشیو

در برآورد هزینه، حاضر فقط هزینه های انجام شده توسط کتابخانه ملی در نظر گرفته شد بدین معنی که هزینه های مربوط به استخدام کارکنان آرشیو که سازمان های شریک در نظر داشتند، حذف گردید.

روش شناسی انجام کار

جهت برآورد هزینه کارکنان اقدام به تهیه نموداری از تمامی فعالیت ها و فرآیندهای (وظایف صورت گرفته توسط بخش آرشیو دیجیتال) مورد نیاز جهت کسب شاخصی در زمینه ایجاد آرشیو کردند. در نمودار مزبور، هزینه های (فعالیت های) اساسی مشخص شده بود، سپس کارکنان اقدام به تخمین مدت زمانی کردند که به طور متوسط هر یک از کارکنان در هر روز صرف می کردند یک روز کاری شامل 441 دقیقه است

• شناسایی و انتخاب مواد و منابع - 30 دقیقه

• تماس با ناشر، مذاکره برای کسب اجازه جهت افزودن منبعی به آرشیو و بایگانی مکاتبات - 30 دقیقه

• گردآوری تضمین کیفیت و آرشیو منابع - 210 دقیقه

• فهرست نویسی - 81 دقیقه

• فعالیت های دیگر (شامل انجام مکاتبات با افراد و نهاد های مسئول در امر نمایه سازی و چکیده نویسی، سوالات مرجع آرشیو مدار و همکاری کارکنان آرشیو دیجیتال برای بسط و توسعه پانداس) - 60 دقیقه

● برقراری ارتباط با شرکا و نیز پشتیبانی - 30 دقیقه (این فعالیت در برآورد هزینه منظور نشده است) البته تمامی کارکنان کلیه وظایف فوق را به یک نسبت انجام نمی دهند برخی بیش تر و برخی کم تر کار می کنند برای مثال سرپرست بخش عهده دار انجام وظایف اداری و نظارت بر کارکنان نیز می باشد این بدان معنی است که او کم تر از دیگران به کارهایی چون گردآوری، آوری کنترل کیفیت، آرشو سازی،

ص: 211

و فهرست نویسی می پردازد.

از میان مجموعه وظایفی که در بالا اشاره شد، مدیر بخش ایجاد آرشیو دیجیتال، تنها متعهد به حمایت از شرکا و برقراری تماس با ناشران است اوقات وی بیش تر صرف، مدیریت اعمال سیاست گذاری های، توسعه تبلیغات و ارتباط با دیگر سازمان های فعال در امر ایجاد آرشیو تحت وب است. با وجود این حقوق و دستمزد وی جزء ضروری هزینه های کلی ایجاد آرشیو دیجیتال کتابخانه محسوب می شود و باید در برآورد کلی هزینه ها منظور شود.

در این میان کارمندانی نیز عهده دار انجام اموری هستند که کاملاً از حیطه وظایف آن ها خارج است. برای مثال، دو نفر از کتابداران هر هفته باید وقت خود را صرف میز مرجع موجود در اتاق مطالعه کند، و باید این مدت زمان را از مجموع مدت زمان حضور آن ها در آرشیو کسر کرد.

کتابخانه ملی به عنوان شریک اصلی پاندورا و تأمین کننده زیر ساخت های فنی دارای نقشی حمایتی نسبت به دیگر شرکاست و این، امر به نوبه خود به صرف هزینه های بیش تری می انجامد. تمامی کارکنان بخش آرشیو دیجیتال مقادیر قابل توجهی از وقت خود را صرف برقراری ارتباط با دیگر همکاران و ارائه خدمات فنی پشتیبانی کننده به آن ها می کند کارمندان، ستادی هر روز به طور متوسط 30 دقیقه از وقت خود را صرف انجام این کار، یعنی برقراری ارتباط با دیگر همکاران و حمایت از آن ها می کند، در حالی که همزمان مدیر آن ها تنها 10 دقیقه از وقت خود را در هر روز صرف انجام این مهم می کند. لازم به ذکر است که هزینه های مربوط به موارد فوق در صورت هزینه های برآورد شده منظور نشده است.

نحوه محاسبه هزینه ها

پس از برآورد هزینه ها در مرحله بعد با استفاده از برنامه اکسل روشی برای محاسبه هزینه های مربوط به هر یک از فعالیت های انجام گرفته در آرشیو، طراحی و تدوین شد.

سپس در هر روز مقادیر مربوط به حقوق و دستمزد کارکنان آرشیو دیجیتال و نیز مقادیر زمانی صرف شده توسط آن ها در خصوص انجام هر کار وارد برنامه صفحه گستر اکسل شد و در نهایت، کل تعداد دقیق صرف شده برای انجام هر کار و نیز هزینه های مربوط به هر کارمند، محاسبه گردید.

سپس هزینه های تأمین کنندگان مختلف به شرحی که پیش تر گفته شد نیز اضافه گردید، و هزینه های مربوط به توسعه زیر ساخت های نگهداری نیز در این مرحله لحاظ شدند.

در مجموع می توان گفت که از کل 937 موردی که توسط کتابخانه ملی در ماه های ژوئیه تا اکتبر سال 2004 وارد آرشیو شد به طور متوسط در هر روز 13 عنوان به ثبت رسید.

هزینه های فراهم آوری منابع آرشیوی

با ورود تمامی این اطلاعات در برنامه صفحه گستر اکسل مشخص شد که هزینه های تمام شده برای هر منبع آرشیو شده، صرف نظر از هزینه های مربوط به انجام فعالیت هایی چون برقراری ارتباط با همکاران و حمایت از آن ها رقمی معادل 178/68 دلار استرالیا است.

اطلاعات مربوط به هر یک از اجزای تشکیل دهنده

ص: 212

هزینه های مربوطه در زیر آمده است:

• هزینه های مربوط به استفاده از نیروی انسانی جهت آرشیو هر منبع در آرشیو دیجیتال - به ازای هر مورد 168/36 دلار استرالیا؛

• هزینه های فراهم کنندگان اقلام آرشیوی به ازای هر مورد 3/41 دلار استرالیا؛ و

• هزینه های مربوط توسعه زیرساخت ها و تعمیر و نگهداری اقلام آرشیوی به ازای هر مورد 6/791 دلار استرالیا

همان گونه که ملاحظه می کنید، حقیقت تلخ پر هزینه بودن این کار سخت و پرمشقت، در رویکرد گزینشی ایجاد آرشیو وب به خوبی نمایان و قابل مشاهده است. هزینه های مربوط به نیروی انسانی، 94 درصد کل هزینه ها را تشکیل می دهد

مقایسه با نوع چاپی

تفاوت هزینه های بالای دستیابی به انتشارات، و بی در مقایسه با نشریات چاپی را همچنین می توان از قیاس میان هزینه های مربوط به تهیه این گونه مواد (وبی) و موارد مشابه چاپی (شامل پیاپند ها و تک نگاشت ها) که به طور همزمان از طریق قانون و اسپاری فراهم گردیده اند، احراز نمود.

• هزینه فراهم آوری تک نگاشت ها با استفاده از قانون و اسپاری منابع چاپی به ازای هر مورد 12/29 دلار استرالیا؛ و

• هزینه فراهم آوری پیایندهای چاپی با استفاده از قانون و اسپاری منابع چاپی به ازای هر مورد 11/29 دلار استرالیا.

تلاش برای مقایسه هزینه های مزبور اندکی شبیه تلاش برای مقایسه هزینه های مربوط به حمل و نقل هندوانه موز و انگور در بازار در حالت عمده در قیاس با قیمت تمام شده برای هر یک از این نوع اقلام است. در حالت عادی قیمت تمام شده هر یک از این اقلام با قیمت عمده آن ها برابر نیست. در قیاس با این مسئله، قیمت و هزینه تمام شده برای تهیه منابع تحت وب، با قیمت تمام شده تهیه هریک از این مواد در حالت چاپی نیز برابری ندارد گر چه اختلاف فاحشی میان ماهیت مواد مزبور (منابع تحت وب) و فرآیند تهیه آن ها با نوع چاپی وجود دارد کاملاً آشکار است که در کل، هزینه تهیه و فراهم آوری منابع تحت وب در مقایسه با منابع چاپی بالا تر می باشد

باید این جا به این نکته توجه داشت که در میان اعداد و ارقام موجود به هیچ عنوان هزینه های مربوط به خرید منابع منظور نشده است. کتابخانه منابع چاپی را بر اساس طرح و اسپاری و بدون پرداخت هزینه ای دریافت می کند منابع تحت وب نیز بدون پرداخت هزینه از طریق وبگاه ناشران قابل دستیابی است. این هزینه ها صرفاً مربوط به فرآیند خرید منابع است هزینه های مربوط به قفسه آرایی و آماده سازی اقلام بر روی قفسه ها، منظور، اما هزینه های بعد آن نظیر و جین منابع و مراقبت از مجموعه در محاسبات مربوط به برآورد هزینه ها لحاظ نمی شود.

افزون بر این کتابخانه ملی علاقه مند به تحلیل هزینه های مربوط به ارائه خدمات خاص نیز بود از این رو جزئیات هزینه های مربوط به استفاده از افراد به ازای انجام هر یک از موارد کاری زیر عبارت است از:

• شناسایی و انتخاب مواد و منابع - 10/16 دلار استرالیا

• تماس با ناشر، مذاکره برای کسب اجازه جهت افزودن منبعی به آرشیو و بایگانی مکاتبات - 10/34 دلار استرالیا

• گردآوری تضمین کیفیت و آرشیو منابع - 7/09 دلار استرالیا،

• فهرست نویسی - 27/42 دلار استرالیا، و

• دیگر فعالیت ها (نظیر بسیاری از فعالیت های مدیریتی انجام مکاتبات با افراد و نهاد های مسئول در امر نمایه سازی و چکیده نویسی پرسش های مرجع آرشیو مدار و همکاری کارکنان آرشیو دیجیتال برای بسط و توسعه پانداس - 59,67 دلار استرالیا.

توجه داشته باشید که این برآورد هزینه چندان نیز مفید نیست زیرا این برآورد هزینه تنها به شکل موردی و در سطح فهرست نویسی ساده و ابتدایی (توصیفی) صورت می پذیرد نه در سطحی تحلیلی و محتوایی در حالی که هر عنوان موجود در آرشیو دارای حداقل دو سطح فهرست نویسی است. به این ترتیب قیمت واقعی ترقمی معادل 54/84 دلار استرالیا می باشد.

امکان کاهش هزینه ها

با اطلاعات حاصل از این روش (تخمین هزینه فعالیت های صورت گرفته) توانستیم دریابیم که چگونه می توان هزینه ها را کاهش داد ما توانستیم هزینه های خود را با ایجاد تغییر در رویکرد و اقدام جهت ایجاد آرشیو تحت وب (تغییر سیاست های کاری خود) با یافتن شیوه هایی جهت انجام مؤثرتر وظایف خود در این زمینه و با اعمال همزمان هر دو روش کاهش دهیم برای دستیابی به آینده ای قابل پیش بینی، تلاش خواهیم کرد تا تدابیر و سیاست های موجود در این عرصه برای مثال تدوین فهرستگان اطلاعات کتاب شناختی منابع یا تضمین کیفیت را حفظ کنیم. البته موقعیت های دیگری نیز برای صرفه جویی در هزینه ها وجود دارد.

شناسایی و انتخاب عناوین، اقدامی مهم و اساسی برای کاربرد رویکرد گزینشی در ایجاد آرشیو دیجیتال است. در نگاه اول و در نخستین گام مشکل به نظر می رسد که در یابیم چگونه می توان مدت زمان صرف شده برای انجام این قبیل فعالیت ها را کاهش داد به هر حال ما با ناشران دولتی از آن جهت همکاری می کنیم که بتوان از طریق آن ها به ابر داده های مربوط به اطلاعات کتاب شناختی منابع برخط که به صورت گروهی توسط آن ها در پانداس بارگذاری شده و امکان دستیابی خودکار به منابع برخط را فراهم می سازد - دست یابیم در واقع ناشران آن چه را که پاندورا آرشیو خواهد شد - به واسطه تأمین ابر داده های مربوط به، آن ها مشخص می نمایند. پانداس همچنان در صدد است تا از این طریق امکان پردازش و دستیابی به منابع اطلاعاتی را به شکلی انبوه برای همگان محقق سازد با وجود این، این امر زمانی محقق

خواهد شد که قیمت متوسط فراهم آوری منابع اطلاعاتی توسط سازمان های همکار کاهش یابد. هم اکنون فراداده ها توسط تعداد اندکی از ناشران به طور خودکار به رکوردهای مارک ارسال می شوند تا ضمن افزوده شدن به پایگاه داده های کتابشناختی، ملی به طور همزمان در فهرستگان برخط کتابخانه نیز بارگذاری شوند این امر به کاهش هزینه ها خواهد انجامید.

محاسبات بیشتر در خصوص جزئیات مربوط به هر فعالیت (جمع آوری، تضمین کیفیت، و آرشیو منابع) بیان گر آن بود که تضمین کیفیت 86 درصد کل هزینه ها را تشکیل می دهد. چنان چه می توانستیم از نرم افزارهای قابل اعتماد در زمینه ارزیابی کیفیت در پانداس استفاده کنیم صرفه جویی در هزینه ها بسیار با ارزش تر و معنی دارتر می شد

افزون بر این کتابخانه مزبور تلاش می کند با انجام واسطه گری و لابی های مناسب، دولت استرالیا را متقاعد به بسط و گسترش امکان استفاده از قانون واسطه گری منابع برخط، کند تا از این طریق بتواند نسبت به رفع موانع موجود در زمینه ضرورت اخذ مجوز از ناشران به منظور تهیه نسخه منابع آرشیوی اقدام نماید. بدیهی است تحقق این امر موجب کاهش زمان کاری کارکنان آرشیو خواهد شد.

بهره گیری از فناوری به پیشرفت این امر کمک خواهد کرد، و ما را قادر خواهد ساخت که هزینه های مربوط به انجام این فعالیت را کاهش دهیم کنسرسیوم بین المللی محافظت از منابع اینترنتی که کتابخانه ملی نیز از اعضای فعال آن می باشد، در صدد تهیه مجموعه ای از ابزار های لازم جهت ایجاد آرشیو وب نظیر میانجی کاربر به طور اخص طراحی شده برای کتابخانه ملی. است انتظار می رود تحقق این امر نیز به کار آمد شدن این فعالیت بیانجامد.

نتیجه گیری

محاسبه هزینه ها، دانسته های کنونی ما را در خصوص ایجاد آرشیو تحت وب و این که این نوع آرشیو در مقایسه با شیوه فراهم آوری منابع چاپی فعالیتی پر هزینه بوده (حتی اخیراً هزینه های آن به دلیل افزایش هزینه های نیروی کار نسبت به گذشته فزونی یافته است) تأیید کرد. علاوه بر این، مطالعه حاضر اطلاعات سودمندتری را درباره چگونگی ارتقای برنامه ها در اختیار ما قرار داد آیا ما باید از سیاست واحدی در زمینه فهرست نویسی عناوین آرشیوی استفاده کنیم؟ بله شاید از این طریق بتوانیم راه های کم زحمت تری را جهت ثبت پیشینه های کتاب شناختی منابع به دست آوریم تا چه میزان باید وقت و انرژی خود را صرف فرآیندگران قیمت و پرهزینه تضمین کیفیت نماییم؟ مسلماً، خیلی اما بهره گیری از یک نرم افزار ارزیابی مناسب می تواند بسیار سودمند باشد.

بدیهی است افزایش پختگی در زمینه بهره گیری از فناوری های مربوط به ایجاد آرشیو تحت وب، منجر به افزایش کارآمدی تمامی پروژه های مربوط به این زمینه (ایجاد آرشیو تحت وب)، نظیر، آرشیوهایی که اخیراً اقدام به استفاده از رویکرد پر زحمت گزینشی کرده و موفق شده اند از هزینه های سرانه پایین تری در یک بازه زمانی برخوردار گردند خواهد شد.

مایلم در اینجا اعلام کنم که بیشترین زحمات در برآورد هزینه‌ها بر عهده ناظم یوسف، مدیر مالی کتابخانه استرالیا بوده است. روش شناسی این کار شامل تخمین هزینه‌ها و تفسیر نتایج مربوط به آن‌ها نیز بر عهده وی بوده است.

منابع

1. Library and Archives Canada. Electronic Collection: a Virtual Collection of Monographs and Periodicals. .1 <http://www.collectionscanada.ca/electroniccollection/> (accessed 22 March 2005); National Diet Library of Japan. Web Archiving Project (WARP). <http://warp.ndl.go.jp/> (accessed 22 March 2005); National Library of Australia. www.nla.gov.au 9|10 Creative Commons Attribution-NonCommercial-ShareAlike 2.1 Australia 19 March 2009 Staff paper Selective Archiving of Web Resources 10 | 10 www.nla.gov.au 19 March 2009 Creative Commons Attribution-NonCommercial-ShareAlike 2.1 Australia PANDORA, (Australia's Web Archive. <http://pandora.nla.gov.au/index.html> (accessed 22 March 2005
2. (National Library of Sweden. Kulturawa3. <http://www.kb.se/kw3/> (accessed 22 March 2005
3. (Library of Congress. MINERVA. <http://www.loc.gov/minerva/> (accessed 22 March 2005
4. State Library of Tasmania. STORS: Long Term Storage of Tasmanian Electronic Documents. . 4 <http://www.stors.tas.gov.au/> (accessed 22 March 2005); Oltmans, E. and H. van Wijngaarden. 2004. Digital preservation in practice: the e-Depot at the Koninklijke Bibliotheek. *Vine* 34 (1): 21-26
5. Information about PANDORA, Australia's Web Archive, and access to its contents are available at . 5 <http://pandora.nla.gov.au/index.html>
6. Information about PANDAS is available at <http://pandora.nla.gov.au/pandas.html> . 6
7. (National Library of Australia. Kinetica. <http://www.nla.gov.au/kinetica/> (accessed 22 March 2005 . 7
8. Levels and pay scales are explained in Attachment A – Salary Table of the . 8

.National Library of Australia Certified Agreement 2004–2007 available here

An "instance" is a single gathering of a title. It includes the gathering of a monograph that has been archived .9
once only, the first gathering of a serial or integrating title (for example, a website that changes over time),
.and all subsequent gatherings

International Internet Preservation Consortium <http://netpreserve.org/about/index.php> (accessed 22 . 10
.March 2005

ص: 217

بایگانی اینترنت اقدامی پیشگامانه و تلاشی ابتکاری برای بایگانی منابع وب انجام داده است. معمولاً جست و جوی این نوع بایگانی ها با استفاده از فراداده های توصیفی یا خدمات نمایه سازی متن، کامل امکان پذیر نیست. مسائل بسیاری وجود دارد که بایگانی وب را مشکل می کند و اگر برنامه ریزی سازمانی، برای بایگانی وب از محدوده کوچک تری آغاز شده باشد و بیش تر محدودیت های آن محدودیت های بشری یا محدودیت های منابع سخت افزاری باشد، به برخی مسائل باید با رویکردی متفاوت نگاه کرد. برای نشان دادن برخی مسائل اصلی، بایگانی رقومی برای مطالعه زبان چینی (SHCAD) به عنوان مطالعه موردی در مقاله حاضر مورد بررسی قرار گرفته است.

هانوالشر (2) | ترجمه حمزه علی نور محمدی (3)

1. جرایب بایگانی علمی در مقیاس کوچک

با توجه به پیچیدگی های بایگانی وب و نیاز های سخت افزاری و نرم افزاری و همچنین تخصص و پرسنل آیا چنین پروژه هایی فقط برای مؤسسه هایی در مقیاس بزرگ مانند کتابخانه های ملی امکان پذیر است؟ یا اینکه مؤسسه های کوچک تر مانند موزه ها دانشگاه ها و مانند آن ها نیز قادر به انجام وظایف مورد نیاز برای بایگانی وب- با چشم انداز بلند مدت قادر به این کار خواهند بود؟

حتی اگر پاسخ مثبت باشد این سؤال باقی می ماند که آیا در واقع این کار ضروری است یا خیر؟ به هر حال می توان فکر کرد که در حال حاضر بایگانی اینترنت همراه با تلاش برای افزایش تعداد کتابخانه های ملی بسیاری از منابع وب را تحت پوشش دارد و بایگانی می کند.

اجازه دهید با سؤال دیگری شروع کنم بایگانی اینترنت اقدامی پیشگامانه را به عنوان نخستین تلاش ابتکاری، جامع جهت بایگانی منابع وب انجام داده است. موفقیت در انجام این کار، انقلابی بوده و بنیاد و اساس پروژه های بسیاری را پی ریزی کرده است با این حال، بررسی آن چه که بایگانی اینترنت و پروژه های جامع دیگر می توانند به آن دست یابند جهت کشف برخی محدودیت ها، آسان است. از آن جا که اساس

ص: 219

Small Scale Academic Web Archiving: DACHS: in Masanes, Julien (ed.), Web Archiving. Berlin. -1

.Heidelberg New York: Springer.pp.213-224

Hanno E. Lecher -2

3- استادیار دانشگاه شاهد

کار مجموعه، بسیار گسترده است، باید برای بخش بزرگی از فعالیت‌ها برای جمع‌آوری اطلاعات وب‌ها به صورت خودکار، تا آن‌جا که امکان پذیر است به روبات‌ها تکیه کرد این نوع جمع‌آوری اطلاعات غالباً بسیار سطحی است بخش‌های زیادی از دست‌رفته بسیاری از صفحات به طور ناقص بارگذاری می‌شوند و بعضی از انواع فایل و نیز وب‌های مخفی به طور کلی نادیده گرفته خواهند شد.

علاوه بر این، از آن‌جا که دریافت به طور خودکار انجام می‌شود در فواصل زمانی نامنظم نمی‌توان انتخابی آگاهانه از منابع داشته باشیم امکان در نظر گرفتن یا تشخیص مطالب مهم موجود - که ممکن است عمر کوتاهی داشته و یا به سختی شناسایی شوند - وجود ندارد، البته بایگانی اینترنت و سایر پروژه‌های بایگانی وب در مقیاس بزرگ برخی امور خاص را انجام داده‌اند که در آن تلاش زیادی صرف توسعه مجموعه‌ها در محدوده خاصی از موضوع‌هایی شده است که انتخاب کرده‌اند. در هر صورت، تعداد این موضوع‌ها بسیار محدود است و طرح‌های پژوهشی زیادی باید آرشیوهای خود را توسعه دهند.

بیش‌ترین مشکل در پروژه‌های جامع بایگانی، وب‌دسترسی محدود به محتوای مهم‌ترین اصل شناسایی نشانی دقیق اینترنتی یک سند یا وبگاه است تا بازیابی اطلاعات مقدور شود. معمولاً جست‌وجوی این نوع بایگانی‌ها با استفاده از فراداده‌های توصیفی یا خدمات نمایه‌سازی متن کامل، امکان پذیر نیست. حتی اگر گزینه جست‌وجوی متن کامل در دسترس باشد عدم انتخاب آگاهانه منابع، تنها نتایج نامنظمی شبیه جست‌وجوی وب امروزی ارائه می‌دهد.

با نگاه به این محدودیت‌ها بدیهی است که بایگانی اینترنت و برخی منابع دیگر، نه بایگانی اینترنت را به طور کامل انجام می‌دهند و نه راه‌های دسترسی مناسبی برای بسیاری از اهداف علمی و یا پژوهشی دیگر فراهم می‌کنند. در نتیجه بایگانی علمی در مقیاس کوچک به یک نیاز مهم تبدیل شده است. اما آیا بازگشت به پرسش مطرح شده در ابتدا امکان پذیر است؟

به برخی مسائل باید با رویکردی متفاوت نگاه کرد اگر برنامه ریزی سازمانی، برای بایگانی وب از محدوده کوچک تری آغاز شده باشد و بیش‌تر محدودیت‌های آن محدودیت‌های بشری/محدودیت‌های منابع سخت‌افزاری باشد برای نشان دادن برخی مسائل اصلی بایگانی دیجیتال برای مطالعه زبان چینی (DACHS) به عنوان مطالعه موردی مورد بررسی قرار خواهد گرفت. هر چند باید در ذهن داشته باشید که پروژه‌های متفاوت غالباً به روش‌ها یا راه‌حل‌های متفاوتی نیاز دارند.

2. بایگانی دیجیتال برای مطالعات زبان چینی

اشاره

اهداف اصلی بایگانی دیجیتال مطالعات زبان چینی (DACHS)، اطمینان از دسترسی درازمدت برای شناسایی و بایگانی منابع اینترنت مربوط به مطالعات زبان چینی است انتخاب نقش مهمی در این فرآیند دارد و به طور ویژه در گفتمان سیاسی و اجتماعی بر آن تأکید و در اینترنت چینی منعکس شده است.

در حال حاضر پروژه DACHS توسط کتابخانه‌های دو مؤسسه چین شناسی با نام‌های «مؤسسه چین» در شهر هایدلبرگ آلمان و «مؤسسه چین شناسی» در دانشگاه لیدن هلند به کار گرفته شده است. بدین ترتیب، زیرساختی کاملاً متفاوت با کتابخانه‌های بزرگ ملی دارد (جدول 1)

جدول 1- حجم ذخیره مجموعه Dchs

نخستین بار هنگامی که اندیشه بارگذاری منابع برخط در موردچین در اواخر 1999 در هایدلبرگ مطرح شد، به هیچ وجه وضعیت آن برای کسی روشن نبود این اندیشه هنوز توسعه نیافته بود و به عنوان بخشی از برنامه ای بزرگ تر برای ایجاد مرکز اروپایی منابع دیجیتال در مطالعات چینی، با هدف بهبود شرایط مربوط به تحقیقات چین و دسترسی به اطلاعات در اروپا، معرفی شد. این پروژه، شامل فعالیت هایی همچون خرید طیف وسیعی از پایگاه داده های تجاری تمام متن حمایت از توسعه پروژه های پایگاه داده دانشگاهی و توسعه امکانات بود، همچنین توسعه برای یافتن بیماری ایدز در منابع چاپی و غیر چاپی در چین و فراهم آوردن امکان دسترسی آزاد به تمام متن منابع به طور گسترده را شامل می شد. این پروژه به مدت پنج سال در نظر گرفته شد و امکانات مالی آن برای بهبود وسایل سخت افزاری موجود و برای استخدام برخی کارکنان - به عنوان دستیاران دانشجویی - صرف شد.

دستور عمل اصلی این پروژه «حداکثر انعطاف پذیری با حداکثر پاسخگویی» بود. این امر به این معنی است که فضایی برای توسعه طرح های تفصیلی بسیاری از پروژه های فرعی و حتی ایده هایی برای پروژه های فرعی جدید مورد نیاز وجود دارد.

در سال 2000، بیشتر به جذب پروژه های فرعی برای اجرای برنامه ریزی پیشرفته و سازماندهی مجدد زیر ساخت های فناوری اطلاعات مؤسسه پرداخته شد، و برنامه ریزی واقعی برای بایگانی وب در سال 2001 آغاز گردید.

2-1. گام های اولیه

با نگاهی به زیر ساخت های موجود که بایگانی وب باید از طریق آن توسعه داده شود، محدودیت ها قابل مشاهده است. بایگانی دیجیتال برای مطالعات زبان چینی توسط کتابخانه مؤسسه به اجرا در آمد. در آن زمان، کتابخانه با استفاده از فناوری اطلاعات با چهار سرور و نزدیک به 100 ایستگاه کاری [جریان کار] نظارت داشت. مسئول نگهداری این زیر ساخت های فناوری اطلاعات کتابداران با کمک یک یا دو دستیار دانشجویی پاره وقت بودند و در صورت نیاز از سوی بخش آی. سی. تی دانشگاه پشتیبانی می شدند.

برای پروژه بایگانی وب باید یک دستیار دیگر دانشجویی به صورت پاره وقت برای تهیه برخی امور (بارگذاری کردن بایگانی ایجاد ابر داده ها و مانند آن) استخدام می شد؛ زیرا کتابدار - علاوه بر مسئولیت های خود برای کتابخانه و محیط فناوری اطلاعات - باید کارهای مدیریت پروژه و توسعه سایر امور نیز در نظر داشته باشد در واقع سهم دستیار پروژه برای توسعه DACHS مهم است. مدرک تحصیلی بالای دستیار برای جلوگیری از تمرکز زدایی دانش در مورد چارچوب نظری پروژه ها اهمیت دارد به خصوص اینکه وقتی کتابدار مؤسسه را جهت کار دیگری محل را ترک کند، اهمیت آن روشن خواهد شد.

البته برخی مسائل، با توجه به اندازه مؤسسه و امکانات آن باید در آغاز پروژه در نظر گرفته شود و به برخی سؤال ها نیز پاسخ داده شود مانند هدف دسترسی دراز مدت منابع بایگانی شده چیست؟ الزامات و نیازهای سخت افزاری و نرم افزاری برای ایجاد و نگهداری بایگانی چه چیزهایی است؟ چگونه انتخاب منابع باید به عنوان یک کار در حال انجام سازماندهی شود و چگونه باید اطلاعات ایجاد شده در دسترس قرار گیرد؟ و مهم تر از همه چه چیز دیگری نیاز است که برای یک برنامه ریزی مناسب در نظر گرفته شود و کجا در جست و جوی پاسخ باشیم؟

پاسخ به آخرین سؤال را می توان در سندی که در سال 2003 به چارچوب استاندارد بایگانی وب تبدیل شده پیدا کرد. 14721:2003 ISO که با عنوان OAIS (سیستم اطلاعات آرشیوی باز) شناخته می شود ثابت کرده است که از اهمیت حیاتی برخوردار است زیرا در آن پس زمینه های تئوری بسیار مهمی تعیین شده است و کمک می کند تا بسیاری از مسائل مهم بایگانی وب را در آن ببینیم.

این سند مفید یک اشکال عمده نیز دارد: همان طور که از نام آن پیداست، فقط چارچوبی برای راهنمایی های تئوری در مسائل مختلف است و پیاده سازی واقعی مهارت های کاربر را ارائه نمی دهد، بنابراین لازم است جست و جوی دیگری برای درک چگونگی تبدیل چارچوب به عمل انجام شود.

اطلاعات بسیاری می توان در جاهایی مانند RLG PADI و مانند آن یافت اما مهم تر از آن مشارکت فعالانه در کارگاه ها و همایش هایی است که در برخورد با مسائل مرتبط با بایگانی وب ارائه می دهند.

2-2. توسعه پایدار سازمانی

یکی از مهم ترین سؤال هایی که کل پروژه با آن روبه روست چگونگی فراهم کردن توسعه پایدار مؤسسه برای بایگانی دراز مدت است. هیچ یک از سه عامل اساسی در این پرسش نمی توانند به صورت دراز مدت در نظر گرفته شوند تأمین مالی پروژه نمی تواند محدود شود و به پایان برسد؛ علاقه مندی های کاربران در سازمان می تواند تغییر کند و منجر به غفلت از پروژه شود؛ و حتی وضعیت دراز مدت سازمان به هیچ وجه تضمین شده نیست. بنابراین واضح است که سازمان آن جایی با قابلیت اعتماد به مراتب کم تری از کتابخانه یا بایگانی ملی برای بایگانی دراز مدت است که در آن برخی مفاد قانونی می تواند مؤسسه را برای تحقق مسئولیت ها نسبت به مجموعه مجبور می کند.

بنابراین، راهبرد توسعه باید برای بقای بایگانی، فنی در صورت توقف کار مؤسسه یا عدم توانایی آن برای حمایت از پروژه باشد بقا را می توان به دوروش تعریف کرد: یا باید بایگانی را فعال نگه داشت بدین معنی که تمام فعالیت های انتخاب منابع برای داده های در دسترس ادامه داشته باشد؛ یا این که بایگانی باید حداقل در یک وضعیت رکودی حفظ شود، بدین معنی که با وجود این که منابع جدید به بایگانی اضافه نمی شوند حداقل دسترسی به آن چه که در حال حاضر موجود است تضمین شود.

دوره برای این کار وجود دارد نخست و مهم تر از همه بایگانی باید ویژگی های اولیه یک منبع قابل اعتمادی را محقق کند، که به معنای پایبندی به استانداردهای اعلام شده است همان گونه که در مدل OAIS توصیف شد. در عمل استفاده از آن را برای مؤسسه های دیگر ممکن سازد - حالت ایده ال کتابخانه های ملی - البته تا زمانی که قادر به همسانی آن با سایر بایگانی ها نباشیم

راه دوم توسعه بایگانی است که برای آن تلاش شده است. اگر تعدادی از مؤسسه ها به طور فعال و تعاملی در این پروژه مشارکت داشته باشند در زمان توقف کار برای یکی از شرکا برای دیگری این امکان به وجود می آید که کار بایگانی را ادامه داده و از امکانات آن استفاده کند. در عمل، پایبندی به استانداردهای برقرار شده در این نوع کار ضروری به نظر می رسد برای این پروژه نیز تصمیم گرفته شد که به صورت مشارکتی کار شود [پروژه] در جست و جوی شرکایی است که امکان همکاری در این زمینه را داشته باشند و مسائل مربوط را حل کنند و [در نهایت] بایگانی به عنوان یک کمک اساسی به حوزه مطالعات چین وارد شده و شناخته شود.

3-2. سخت افزار

برای اطمینان از کارکرد مناسب DACHS در سطح محلی باید محیط سخت افزاری مناسبی را اندازه گیری شود بعد از تدارک سخت افزاری برای سرور مناسب و ایستگاه های کاری اختصاص داده شده به امور روزمره مانند بارگذاری و مقاصد، مدیریتی لازم است به امکانات پشتیبانی و امنیت و محافظت از این سرور در مقابل ویروس های کامپیوتری توجه شود.

مرکز کامپیوتر دانشگاه هایدلبرگ با استفاده از روش ذخیره سازی توزیع شده IBM ADSTAR، سیستم پشتیبان گیری مهمی را فراهم کرده است. با استفاده از این سیستم هر شب یک نسخه پشتیبان از کل بایگانی بر روی نوار های مغناطیسی در کامپیوتر مرکزی ذخیره می شود کپی های پشتیبانی از این نوارها نیز به طور مرتب در دانشگاه کالسروهه ذخیره می شود بنابراین داده های بایگانی در سه مکان مختلف که امنیت مناسبی را ارائه می دهند نگهداری می شوند.

دستگاه منبع تغذیه اضطراری (UPS) نیز به عنوان یک سیستم RAID (آرایه افزونه دیسک های مستقل ارزان) (سطح اول) جهت تأمین امنیت اولیه برای دسترسی بی وقفه نصب شد. برای محافظت در مقابل ویروس از نرم افزار ویروس یاب McAfee استفاده شد. از طریق تعاریفی که در این ویروس یاب وجود دارد، به صورت ساعتی از سرور McAfee استفاده می شود تمام داده های ورودی به طور مداوم چک شده و به طور خودکار فرآیندهای اسکن منظم از کل بایگانی انجام می پذیرد.

بدینوسیله می‌توانیم زیر ساخت‌های فناوری اطلاعات موجود در مؤسسه را که حمایت بخش ICT دانشگاه را نیز دارد در این بخش مورد استفاده قرار دهیم در تمامی موارد این تجهیزات باید از نظر دانشگاه‌ها استاندارد، باشد به همین دلیل این تجهیزات باید در حوزه آی.سی.تی. مرکزی بوده و در آن جا نگهداری شود. این موضوع می‌تواند از لحاظ حرفه‌ای نوعی مزیت به حساب آید. از طرفی هم می‌تواند یک نقص باشد زیرا به نوعی می‌تواند محدودیت‌هایی بر روی سخت‌افزار و نرم‌افزار اعمال کند.

4-2. نرم افزار

برای شروع کار بایگانی دیجیتال برای پروژه مطالعات چینی به نرم افزار مناسب نیاز داریم که باید برخی شرایط این کار را داشته باشد.

اینترنت مجموعه عظیمی از داده‌های به هم پیوسته در فرمت‌ها و سیستم‌های کدگذاری مختلف است و بایگانی چنین داده‌هایی برای محافظت از آن‌ها به نحوی که محتوا و عملکرد آن به صورت اصلی باشد، قدری مشکل به نظر می‌رسد. از آن جا که کارکرد محتوای وب وابسته به نرم افزار مرورگر است، هیچ داده قابل اعتمادی برای محافظت از آن وجود ندارد با این حال چنان چه اساس آن محافظت شود، بسیاری از موارد فوق برآورده می‌شود به این معنی که اگر کل وبگاه بارگذاری شد، باید ساختار فایل اصلی دست نخورده نگه داشته شود البته برای نگهداری اسناد بارگذاری شده و ارتباط آن با باید ارتباطات مناسبی در نظر گرفته شود برای این منظور به نرم افزاری نیاز است که مناسب این کار بوده و قابلیت تنظیم را داشته باشد این نرم افزار باید توانایی اداره طیف گسترده‌ای از فرمت‌های مختلف را داشته و مقرون به صرفه نیز باشد بعد از آزمایش‌هایی که بر روی نرم افزارها صورت گرفت به این نتیجه رسیدیم که نرم افزار آفلاین اکسپلورر، تمام قابلیت‌های مورد نیاز را برای شروع کار دارد.

5-2. فراداده

مسئله‌ای که موجب دغدغه فکری شده ایجاد فراداده است از یک طرف فراداده اطلاعات لازم را برای مشاهده و توصیف محتوای سند ایجاد می‌کند و از طرفی هم موجب دسترسی کاربران به اطلاعات می‌شود. فراداده‌هایی مانند نویسنده عنوان و موضوع تا زمانی که جست‌وجوی متنی وجود نداشته باشد این نوع فراداده به کاربر کمک می‌کند از طرفی هم فراداده مسائل فنی را برای کمک به حفاظت درازمدت فراهم می‌کند؛ زیرا اطلاعات لازم را برای مدیریت مناسب را فراهم و بررسی را آن در آینده میسر می‌کند.

ایجاد فراداده پر هزینه است؛ با این حال با وجود فراداده نیمه خودکار، جمع‌آوری برخی داده‌ها باید به صورت دستی انجام شود. بنابراین بدیهی است که ایجاد فراداده برای صدها هزار سند با همان سرعتی که از اینترنت بارگذاری شده غیر ممکن است (شکل 1 را ببینید). حتی اگر این کار ممکن بود یافتن آن چه که دقیقاً این فراداده باید حاوی آن باشد، باز هم یافتن جزئیات و فرمت‌های آن، بسیار سخت است.

یکی از پرسش‌های مهمی که در این مورد وجود دارد این است که آیا وجود فراداده ضروری است یا بهتر است بگوییم: آیا صرفاً تکیه بر الگوریتم‌های جست‌وجو برای بازیابی تمام داده‌های مورد نیاز

برای دسترسی و نیز برای اهداف بلند مدت حافظت از داده‌ها امکان پذیر است؟ از دیدگاه کاربر می‌توان استدلال کرد که جست و جوی متن، کامل در مقایسه با سرفصل‌های موضوعی خام یا عناوین حاوی فراداده‌های ناقص جهت پیدا کردن اسناد ابزار قابل اعتماد تری است از نظر فنی بسیاری از اطلاعات ضروری از جمله فرمت فایل، سرعت بارگذاری، رمز گذاری و غیره برای راحتی در بازیابی داده‌هاست یا حتی می‌تواند برای ساختار نام گذاری فایل باشد

عکس

بایگانی وب علمی در مقیاس کوچک DACHS ۲۲۵

برای دسترسی و نیز برای اهداف بلند مدت حافظت از داده‌ها امکان پذیر است؟ از دیدگاه کاربر می‌توان استدلال کرد که جست و جوی متن کامل، در مقایسه با سرفصل‌های موضوعی خام یا عناوین حاوی فراداده‌های ناقص جهت پیدا کردن اسناد ابزار قابل اعتماد تری است. از نظر فنی، بسیاری از اطلاعات ضروری از جمله فرمت فایل، سرعت بارگذاری، رمز گذاری، و غیره برای راحتی در بازیابی داده‌هاست یا حتی می‌تواند برای ساختار نام گذاری فایل باشد.



شکل ۱-۱. رابط کاربری جست و جوی فراداده DACHS که امکان جست و جوی پیشرفته در بایگانی را دارد.

در پایان این ایده رد شد. منابع دیجیتال نوعی اطلاعات هستند که زود از بین می‌روند. بنابراین امکان دارد تا در اسناد جداگانه‌ای نیز ذخیره شوند تا این منابع را با اطمینان و به صورت یکپارچه در آینده و به صورت طولانی مدت قابل دسترسی کرد. با توجه به نوع داده‌ها و اطلاعات، امکان نبودن آنها وجود دارد. برای مدیریت و ادغام بهتر مواد بایگانی شده در موجودی کتابخانه‌های بزرگ، تصمیم گرفته شد تا فراداده به‌عنوان بخشی از کاتالوگ به صورت منظم ایجاد و تهیه شود. بنابراین، کاتالوگ با فراداده تطبیق داده شد تا بخش‌هایی برای مدیریت درست، تاریخچه مبدأ، سابقه مدیریت، نوع فایل، شناسه، و سایر موارد تعیین شود. بسته به پیچیدگی منابع، در حال حاضر سرگرم تهیه رکوردهای فراداده‌ای هستیم که یا بتواند تک‌فایل را توصیف کند که فقط متن را شامل می‌شود یا شامل مجموعه‌ای کامل از فایل‌ها مانند وبگاه، انجمن و یا روزنامه باشد.

۲-۶. سیاست گذاری و خط‌مشی مجموعه

کاملاً بدیهی است که هدف از بایگانی دیجیتال مطالعات چینی نمی‌تواند شامل تمامی موارد موجود در اینترنت باشد و از آن محافظت نماید. نه از لحاظ فنی چنین امکانی وجود دارد و نه مفید به‌نظر می‌رسد. به‌عنوان مؤسسه‌ای تحقیقاتی، علاقه‌مند به بخشی از اینترنت چین هستیم که منعکس کننده جنبه‌های

شکل 101 رابط کاربری جست و جوی فراداده DACHS که امکان جست و جوی پیشرفته در بایگانی را دارد.

در پایان این ایده رد شد منابع دیجیتال نوعی اطلاعات هستند که زود از بین می روند. بنابراین امکان دارد تا در اسناد جداگانه ای نیز ذخیره شوند تا این منابع را با اطمینان و به صورت یکپارچه در آینده و به صورت طولانی مدت قابل دسترسی کرد. با توجه به نوع داده ها و اطلاعات امکان نابودی آن ها وجود دارد.

برای مدیریت و ادغام بهتر مواد بایگانی شده در موجودی کتابخانه های بزرگ، تصمیم گرفته شد تا فراداده به عنوان بخشی از کاتالوگ به صورت منظم ایجاد و تهیه شود. بنابراین کاتالوگ با فراداده تطبیق داده شد تا بخش هایی برای مدیریت درست، تاریخچه مبدأ، سابقه مدیریت، نوع فایل، شناسه، و سایر موارد تعیین شود. بسته به پیچیدگی، منابع در حال حاضر سرگرم تهیه رکوردهای فراداده ای هستیم که با بتواند تک فایل را توصیف کند که فقط متن را شامل می شود یا شامل مجموعه ای کامل از فایل ها مانند، وبگاه انجمن و یا روزنامه باشد.

2-6. سیاست گذاری و خط مشی مجموعه

کاملاً بدیهی است که هدف از بایگانی دیجیتال مطالعات چینی نمی تواند شامل تمامی موارد موجود در اینترنت باشد و از آن محافظت نماید نه از لحاظ فنی چنین امکانی وجود دارد و نه مفید به نظر می رسد. به عنوان مؤسسه ای، تحقیقاتی علاقه مند به بخشی از اینترنت چین هستیم که منعکس کننده جنبه های

ص: 225

خاصی از جامعه چینی است که موارد زودگذری هم هستند بنابراین فراهم کردن انتخاب آگاهانه منابع یک سرمایه است که در حال حاضر به ما کمک می کند و نیز به کاربران کمک می کند در آینده تا جهت شناسایی موادی که ما نیاز به آن داریم و مربوط به هدف ما می شود یاری می رساند. البته ارزش ها در گذر زمان تغییر خواهد کرد و نسل بعد ممکن است انتخاب های مختلفی داشته باشد. در این مورد، بایگانی اینترنت هنوز جایگزینی بسیار غنی را فراهم میکند که در آن گزینه های دیگری در دسترس خواهد بود و ترکیبی از دو رویکرد گزینشی و جامع است که کاربر در آینده با گسترده ترین آرایه از امکانات برای تحقیقات خود قادر به تأمین آن است.

با این حال همزمان میخواهیم خطر محدودیت بیش از حد را در انتخابمان کاهش دهیم. همانطور که در قبلا نیز گفته شد و در ادامه بحث خواهد شد DACHS در حال تبدیل شدن به پروژه تعاونی بزرگ تر است شرکای، مختلف معیارهای مختلفی دارند و در نتیجه مقادیر مختلفی را در روند انتخاب خود اعمال می کنند این امر تنها می تواند محتوای بایگانی را غنی، کند در حالی که در سیاست انتخاب آگاهانه منابع تغییری رخ نمی دهد

با این حال به منظور استفاده بهتر از منابع محدود و در دسترس باید راهبردهای هوشمندانه ای را برای توسعه مجموعه به کار گرفت. با توجه به تعداد زیاد سیاست های اینترنتی در مورد چین و سرعت زیاد آن و توسعه آن و همچنین ناپدید شدن این مطالب در شبکه چین وظیفه ساختن بایگانی منابع برخط که بتواند مهم و ارزشمند باشد کاری سخت و دلهره آور است.

برای رفع این مشکل شروع به ساخت یک شبکه اطلاعاتی از افراد (محققان بومی و شهروندان) کردیم که به طور فعال یا غیر فعال بخشی از این کار را بر عهده گیرند. این کار زمانی عملی می شود که با استفاده از دانشی باشد که برای شناسایی سیاست های خاص و متناسب با کار باشد با این حال این شبکه اطلاعاتی بزرگ می تواند خیلی متنوع تر از فرآیند انتخاب باشد.

معمولاً شبکه اطلاع رسانی نیز هدف دیگری را دنبال می کند. از آن جا که اعضای آن بخشی از فرهنگ اینترنت معاصر است آن ها می توانند اطلاعات زمینه ارزشمندی در مورد منابع بایگانی شده ارائه دهند تا آن جا که ممکن است این اطلاعات باید در فراداده و یا صفحات وب تخصیص داده شده به عنوان بخشی از مجموعه ایجاد و حفظ گردد.

از آن جا که چند ماه نیز بر روی تعدادی پروژه خاص مانند شعر معاصر، سارس (1) یا موضوع همجنس بازی در چین کار کرده ایم، محققان و دانشجویان کارشناس ارشد برای ایجاد آرشیوهای جامع بر روی این موضوع ها کار می کنند، از جمله این وب گاه ها، آن هایی هستند که در معرض خطر بوده و یا شامل موارد دیگر مانند عکس و پوستر می شود برای بایگانی این منابع می توان متن ساده ای را در زمان بایگانی به آن اضافه کرد این نوع عملکرد را برای درک بهتر این منابع که زود از بین می روند، به ویژه در آینده دور، ضروری می دانیم.

اما این کار صرفاً با تکیه بر آگاهان و دانشمندان انجام نمی شود. تأثیر رویدادهای خاص بین المللی مانند حمله تروریستی 11 سپتامبر و یا بازی های المپیک چین در سال 2008، غالباً باعث بحث های داغ بر

روی اینترنت می شود برای جلوگیری از شیوع افکار عمومی ما در حال کار بر روی سیاهه واریسی تالار گفتگو و روزنامه ها هستیم که نوع گزارش ها به آن مربوط می شود تا بتوان فاصله زمانی چند هفته قبل و بعد از چنین رویدادهایی را پوشش دهد.

مجموعه های مشابه دیگری هم توسط اشخاص، خصوصی، پژوهشگران گروه های پژوهشی و یا سایر مؤسسه های به DACHS اهدا و یا فروخته شده اند. گاهی اوقات برای کمک به ناشران وب گاه ها در معرض خطر برای حفظ محتوای وبگاه شان نزدیک شدیم (و یا با آن ها تماس گرفتیم). به هر حال، این مجموعه ها برای DACHS طراحی و ایجاد نشده اند (یا حتی ممکن است همه استانداردهای کیفیت مد نظر ما را دنبال نکنند) اما DACHS به آن ها کمک می کند تا این منابع در درازمدت در دسترس باشند.

2-7. مشارکت

ملاحظه کردیم که اجرای بایگانی تلاشی مشارکتی است و مشارکت راهبردی مهمی برای بقای آرشیوهای وب در مقیاس کوچک است مشارکت می تواند به نگهداری طولانی مدت از طریق عرضه خدمات به اعضا کمک کند. بنابراین، شرکا نیز اجازه توزیع کار را می دهند که این کار به معنای کاهش هزینه و امکان انتخاب گسترده تر از منابع بایگانی است استانداردهای سخت افزار تجربه و کیفیت نیز می تواند به طور وسیعی به اشتراک گذاشته شود و عملکرد کلی بایگانی را بهبود بخشد. مشارکت به عنوان بخشی از راهبرد سیاسی مهم است. پروژه بین المللی مشارکتی قطعاً ساده تر از شناخت جامعه علمی و کمپین موفق برای تأمین بودجه است.

شناخت مسائل DACHS چیزی است که توسعه دستور عمل ها برای مشارکت ها را ممکن می نماید این دستور عمل ها توسط دانشگاه لیدن، هلند مورد بررسی قرار گرفت که در انتهای سال 2003 مشارکت آن عملی شد. هدف اصلی این بود که شرکا مستقل مانده و تا حد امکان هویت خود را حفظ کنند ولی برخی استانداردها و خدمات وجود دارد که باید به اشتراک گذاشته شود.

یکی از مسائل اصلی در همکاری ایجاد یافته های ایدز، است از جمله راهنمای موضوعی پیونددار (یا جدول محتوا) و همچنین متن کامل و یکپارچه سازی امکانات جست و جو در فراداده و بایگانی آن مسائلی است که همکاری شرکای فعلی و آینده پروژه را می طلبد این امر نه تنها به دسترسی متمرکز به فرم صفحه اصلی به صورت اشتراکی باید دارد بلکه باید در مورد چگونگی گزینه های جست و جو برای رسیدن به هدف مطلوب نیز مذاکره شود جست و جوی متن کامل باید در تمام حوزه ها امکان پذیر باشد و برای فراداده یا یک فهرستگان لازم است که همه داده های فیزیکی با هم ذخیره شود که بتواند از طریق فراداده ها منابع مختلف حتی منابع محلی را جست و جو کند و آن ها را به شیوه ای منسجم ارائه دهد. در هر مورد به ایجاد استانداردهای مختلف جهت ایجاد فراداده نیاز داریم.

مسائل دیگری که باید مورد بحث قرار گیرند عبارت اند از سیاست محدود کردن دسترسی مشترک، مقررات تقسیم کار و امور روزمره ای که از دوباره کاری منابع بایگانی جلوگیری می کند.

نکته مهم دیگر برای همکاری این است که کدام یک از شرکا در آینده قادر هستند زیر ساخت بایگانی به صورت محلی ایجاد نمایند و یا اینکه از امکانات خود برای این پروژه استفاده کنند. از آن جا که

تمامی کارهای فوق هنوز در شرف انجام می باشد به اشتراک گذاشتن تجربه در این زمان غیر ممکن است.

بیش تر مسائل در فرآینده مذاکره و پیاده سازی واقعی برای امور مشارکتی پدیدار خواهد شد. در هر صورت، مشارکت برای پروژه های بایگانی وب در مقیاس کوچک در بسیاری از سطوح ضروری است و از امور مهم و بدیهی است.

3. درس های آموخته شده: جمع بندی

در حدود 4 سال پس از اجرای DACHS، در حال حاضر بسیاری از پرسش ها و مشکلات هنوز حل نشده باقی مانده و یا در حال برطرف شدن و در حال توسعه است. با نگاهی به گذشته، می توان گفت که مسائل کمی وجود دارد که لاینحل باقی مانده است که سعی خواهیم کرد آن ها به طور متفاوتی نگاه کنیم.

ضرورت امر در این زمینه تخصیص سمت ها برای پروژه است. علاوه بر روال کار روزانه باید به کار توسعه و مدیریت پروژه های بایگانی وب نیز توجه کرد و آن را دست کم نگرفت. تصمیم به انحصاری کردن آن و سپردن همه چیز به یک کتابدار که در حال حاضر مجموعه ای از وظایف دیگر را انجام می دهد، قدری جای تأمل دارد.

در هر صورت نقش او در توسعه پروژه مزایای بسیاری دارد و مهم است سمت مدیریتی تخصیص داده شده برای DACHS بسیار مناسب بوده و منجر به اجرای موفق پروژه می شود. جا داشت تا بسیاری از مسائل خیلی زود تر و مؤثر تر به کار گرفته میشد تا تأثیر مثبتی در توسعه پروژه داشته باشد.

مسئله دوم که - حداقل از امروز - باید به طور ویژه ای به پرداخته شود و مسائل آن حل شود انتخاب نرم افزار جمع آوری داده و بایگانی می باشد. در سال 2001 که DACHS شروع به بایگانی وب کرد هنوز در مرحله ابتدایی بود و بسیاری از دست اندکاران بزرگ امروز تنها فقط شروع به توسعه و انتشار تلاش های خود کردند. اگر چه نرم افزار امروز ما که برای اهداف مورد نظر در آن زمان انتخاب کردیم بیش تر نیازهای ما را برآورده می کند ابزار های امروزی نیز برای این کار بسیار مناسب هستند. لازم است ذکر گردد که یک انتخاب جدید باید روند استفاده و ایجاد فراداده را ساده کند.

مطالبی که عنوان شد نقطه نهایی بوده است که من می خواستم آن را ایجاد کنم. مسائل بسیاری هم وجود دارد که بایگانی وب را مشکل می کند و در برخی موارد به مسائلی بر می خورید که به دفعات سراغ آن رفتید ولی خودتان نتوانستید آن را حل کنید بدون اینکه نگران نوع فرمت فایل یا توسعه نرم افزار جمع آوری داده باشیم این امر امکان پذیر است - حتی ضروری است - که بر تلاش دیگران تکیه کنیم. همکاری نه تنها در میان شرکای پروژه بایگانی است بلکه همکاری در سایر مؤسسه های در زمینه بایگانی وب، می تواند راه حل های مهم و ابزارهایی را فراهم کنند که شما به تنهایی قادر به ایجاد آن نخواهید بود.

4. منابع مفید

فهرست زیر تنها انتخاب بسیار کوچکی از منابع است که برای کار مفید است اگر شما وب سایت های زیر را مشاهده کنید منابع بسیار بیش تری را خواهید یافت به خصوص این که این موضوع بسیار عالی با

عنوان دروازه دیجیتالی PADI. نام دارد که همه منابع در ماه می 2006 قابل دسترس بود.

Websites 4.1

Council on Library and Information Resources (CLIR)

<http://www.clir.org/>

Electronic Resource Preservation and Access Network (erpaNet)

<http://www.erpanet.org/>

International Internet Preservation Consortium

<http://netpreserve.org/>

Internet Archive

<http://www.archive.org/>

Networked European Deposit Library (NedLib)

<http://www.kb.nl/coop/nedlib/>

PADI Preserving Access to Digital Information (National Library of Australia)

<http://www.nla.gov.au/padi>

PADI: Web archiving

<http://www.nla.gov.au/padi/topics/92.html>

Mailing Lists 4.2

Archivists

<http://groups.yahoo.com/group/archivists/>

DigiCULT

<http://www.digicult.info/pages/subscribe.php>

DIGLIB - Digital Libraries Research mailing list (IFLA)

<http://infoserv.inist.fr/wwsympa.fcgi/info/diglib/>

OAIS Implementers (RLG)

<http://lists2.rlg.org/cgi-bin/lyris.pl?enter=ois-implementers>

PadiForum (National Library of Australia) <http://www.nla.gov.au/padi/forum/> Web-Archive/

<http://listes.cru.fr/wws/info/web-archive>

Newsletters and Magazines 4.3

CLIR Issues

<http://www.clir.org/pubs/issues/> DigiCULT Newsletter

<http://www.digicult.info/pages/newsletter.php>

D-Lib Magazine

<http://www.dlib.org/DPC/PADI> What is new in digital preservation

<http://www.nla.gov.au/padi/qdiges>

RLG DigiNews <http://www.rlg.org/en/page.php?Page-ID=12081>

ص: 229

هدف این پژوهش بررسی قابلیت های قالب های یونی مارک و مارک 21 برای سازماندهی منابع اطلاعاتی وب و مقایسه آن ها با یکدیگر است تا از این طریق بتوان به ارزیابی این دو قالب در تهیه پیشینه های اطلاعاتی برای منابع اطلاعاتی وب دست یافت روش پژوهش تحلیل محتوا، است به صورتی که ابتدا به شناسایی عناصر توصیف مرتبط با منابع اطلاعاتی وب در قالب کتاب شناختی یونی مارک و مارک 21، و سپس به تحلیل و مقایسه آن ها با یکدیگر پرداخته شد. از یافته های این پژوهش می توان چنین دریافت که قالب مارک می تواند به عنوان استاندارد برای ذخیره و بازیابی منابع اطلاعاتی وب به کار رود با این حال مارک راه حلی برای چگونگی روزآمد نمودن و تعیین سطح توصیف منابع اطلاعاتی وب تهیه نکرده است. کلید واژه ها منابع اطلاعاتی، وب، یونی مارک مارک، 21 قالب کتاب شناختی، سازماندهی

رقیه حجازی، (1) دکتر مرتضی کوکبی (2)

1- مقدمه و بیان مسئله

منابع اطلاعاتی وب بر اساس ویراست 2005 قواعد فهرست نویسی انگلومریکن، شامل منابع (داده ها و یا برنامه) کد گذاری شده به منظور کاربرد با دستگاه های رایانه ای است که نیاز به برقراری ارتباط با یک شبکه رایانه ای دارد (Weitz 2006) این منابع محسوس و شامل محمل فیزیکی نیستند، و همچنین برای استفاده از آن ها نیاز به برقراری ارتباط با یک دستگاه رایانه (مانند شبکه) یا منابع ذخیره شده بر روی دیسک سخت یا دیگر ابزارهای ذخیره سازی است (Miller 2008).

با ظهور محمل های جدید، اطلاعات از جمله منابع اطلاعات موجود در محیط وب نظام های سازماندهی اطلاعات نیازمند روش های نوینی در ذخیره و بازیابی اطلاعات شده اند. نمایه سازی موتورهای کاوش را نمی توان روش مناسبی در ذخیره و بازیابی اطلاعات دانست زیرا با توجه به گفته فتاحی (1380) آشفتگی اینترنت در حال حاضر بیشتر ناشی از پراکندگی و ناکارآمدی روش های سازماندهی اطلاعات در این محیط است بنابراین برای سازماندهی اطلاعات نیاز به رویکردهایی استاندارد با توجه به ویژگی های محمل های اطلاعاتی است

ص: 231

1- کارشناس ارشد علوم کتابداری و اطلاع رسانی 2 r.hejazi86@yahoo.com

2- استاد کتابداری و اطلاع رسانی دانشگاه شهید چمران اهواز kokabi80@yahoo.com

نظام سازماندهی دانش برای سازگاری با تغییر ساختار منابع اطلاعاتی همواره دورویکرد را برگزیده است. نخست هماهنگی سازگاری و تطابق نظام های سنتی با محیط و رسانه های جدید و دیگر طراحی و ایجاد نظام های جدید به منظور حداکثر بهره وری از امکانات و قابلیت های محیط جدید (طاهری 1387). با توجه به اینکه در محیط های کتابخانه ای حجم عظیمی از منابع اطلاعاتی با استفاده از استانداردهای موجود سازماندهی شده اند و تبدیل و تغییر در آن ها نیازمند صرف وقت و هزینه است، به نظر می رسد استفاده از رویکرد نخست حداقل در محیط های کتابخانه ای مناسب تر باشد یکی از استانداردهای مورد استفاده در نظام های سنتی قالب مارک است قالب مارک از استانداردهای مورد استفاده برای ماشین خوان کردن داده های کتابشناختی است. قالب های مختلفی از مارک همچون مارک 21 و یونی مارک وجود دارند. از سوی دیگر، بیانیه کتابخانه کنگره برای چارچوب کتابشناختی برای دوران دیجیتال (2011 در عمرانی 1390) به این نکته اشاره می کند که «استاندارد مارک مسئولیت تولید میلیون ها پیشینه کتاب شناختی در گوشه و کنار دنیا را بر عهده دارد و نیاز است که در تمام دوره نقل و انتقال به پشتیبانی از مارک ادامه دهیم. علاوه بر این بیانیه عنوان می کند که سامانه ها و خدمات مبتنی بر پیشینه های مرتبط با قالب مارک تا سال های زیادی بخش مهمی از زیر ساخت ها خواهند بود». بنابراین توجه به این قالب در سازماندهی اطلاعات حائز اهمیت است زیرا علاوه بر استفاده کنونی از این، قالب الگوها و استانداردهای جایگزین آینده نیز سازگار با آن طراحی خواهند شد.

همان طور که بیان شد قالب های مختلفی از مارک موجود است که می توان قالب مارک 21 (به دلیل کاربرد آن در پیشینه های کتاب شناختی کتابخانه کنگره آمریکا، کتابخانه بریتانیا، کتابخانه ملی کانادا و برخی کتابخانه های دیگر) و قالب یونی مارک (به دلیل بین المللی بودن و مبنای مارک ایران) را از مهم ترین آن ها دانست. هر دو قالب بر اساس قواعد فهرست نویسی انگلومریکن و تأکید بر منابع چاپی تهیه شده اند. با ظهور محمل های اطلاعاتی، جدید تغییراتی در قالب ها ایجاد شد تا بتوانند با محمل های جدید سازگار و در جهت سازماندهی آن ها کارآمد باشند که از آن جمله می توان به افزودن فیلدها و فیلدهای فرعی، نشان گرها و نویسه های موجود در برجسب پیشینه اشاره کرد.

با توجه به ویژگی منابع اطلاعاتی وب از یک سو و تفاوت های موجود میان قالب مارک 21 و یونی مارک این پرسش پیش می آید که آیا قالب هایی همچون یونیمارک و مارک 21 توانایی سازماندهی منابع اطلاعاتی وب را دارند؟ علاوه بر این تفاوت های موجود در میان این دو قالب تأثیری بر کیفیت سازماندهی منابع اطلاعاتی وب دارد؟ پژوهش حاضر سعی دارد با بررسی دو قالب یونی مارک و مارک 21، به بررسی عناصر تعریف شده بر اساس خصوصیات منابع اطلاعاتی، وب و همچنین تبیین کارایی و قابلیت دو قالب در سازماندهی بهینه منابع اطلاعاتی وب بپردازد. بنابراین به ارزیابی توان و کارایی هر دو قالب در سازماندهی اطلاعات وب و مقایسه آن ها با یکدیگر پرداخته شد.

2- پرسش های پژوهش

این پژوهش در پی پاسخ به پرسش های زیر انجام شد:

1-1 هر یک از قالب های یونیمارک و مارک 21 شامل چه عناصری (اعم از فیلد فیلد، فرعی نشان گر و غیره) برای سازماندهی منابع اطلاعاتی وب هستند؟

1-2 آیا عناصر تعریف شده در دو قالب یونی مارک و مارک 21 نشان دهنده خصوصیات منابع اطلاعاتی وب هستند؟

1-3 آیا تفاوتی میان عناصر تعریف شده در قالب یونی مارک و مارک 21 برای سازماندهی منابع اطلاعاتی وب وجود دارد؟

1-4 اگر پاسخ پرسش سوم مثبت است تفاوت های موجود چه تأثیری بر سازماندهی منابع اطلاعاتی وب می گذارد؟

3-پیشینه پژوهش

بررسی پیشینه های مرتبط با پژوهش نشان داد که پژوهش های صورت گرفته بر اساس دو موضوع کلی قابل تقسیم بندی هستند. موضوع نخست در رابطه با ماهیت منابع اطلاعاتی وب و شناسایی تفاوت های آن ها با دیگر قالب های اطلاعات بود که از آن جمله می توان به پژوهش وارد (2001) (1) اشاره نمود. وی در این پژوهش به اهمیت سازماندهی منابع اینترنتی اشاره کرده و معتقد است که مهارت های فهرست نویسی برای سازماندهی منابع موجود بر روی وب مناسب است در این مقاله به برخی از فیلدهای مارک 21 که در سازماندهی منابع اینترنتی کارآمد هستند نیز اشاره شده است. همچنین حاجی زین العابدینی (1381) به بررسی و تحلیل آخرین فعالیت ها و تحقیقات در زمینه سازماندهی اطلاعات در اینترنت و ارائه یک الگوی مناسب برای فهرست نویسی منابع فارسی پرداخت. نتایج به دست آمده نشان داد که دوروش مهم برای سازماندهی اطلاعات در اینترنت وجود دارد. روش نخست ایجاد پیشینه های کتاب شناختی با استفاده از قواعد عام فهرست نویسی و قواعد خاص منابع اینترنتی و روش دوم ایجاد پیشینه های کتاب شناختی با استفاده از روش های ابر داده عنوان شده است. او به دلیل ویژگی های منابع اینترنتی فارسی روش دوم را برای منابع اینترنتی فارسی مناسب ندانسته است.

موضوع دیگر در رابطه با بررسی انواع فراداده ها در سازماندهی منابع اطلاعاتی وب بود که از آن جمله می توان به پژوهش فارد و ریگیو (2004) (2) اشاره نمود که به بررسی چند استاندارد فراداده از جمله قالب مارک برای مدیریت منابع الکترونیکی (که منابع اطلاعاتی وب زیر مجموعه ای از آن است) پرداختند تحلیل های این پژوهش گران مشخص می کند که در زمان انجام پژوهش هیچ استاندارد و الگوی فراداده ای توانایی نمایش پیچیدگی های منابع الکترونیکی را نداشته است. طاهری (1387) نیز به مقایسه کارایی طرح فراداده ای هسته دوپلین و قالب فراداده ای مارک 21 در سازماندهی منابع اطلاعاتی وب پرداخت نتایج پژوهش وی نشان داد که قالب مارک 21 برای ذخیره پردازش و مبادله اطلاعات محیط وب مناسب تر است. علاوه بر این پژوهش های دیگری همچون مرور قالب های ابر داده (Heer)

ص: 233

Ward -1

Fard and Riggio -2

1996)، کنترل مستند در زمینه کنترل کتاب شناختی در محیط الکترونیک (Gorman 2004)، بررسی و مقایسه قواعد فهرست نویسی انگلومریکن و عناصر هسته دویلین برای سازماندهی منابع اینترنتی (کوکبی و آخشیک 1385 و سازماندهی صفحات وب با استفاده از نظام های رده بندی کتابخانه (ملای مقدم و نعیم آبادی 1385) نیز با هدف بررسی و مقایسه دو رویکرد سنتی و نوین فهرست نویسی منابع اینترنتی صورت گرفته است

مطالعه پژوهش های مرتبط نشان داد که سازماندهی منابع وب با استفاده از استانداردها و قواعد فهرست نویسی مورد توجه بوده است. علاوه بر این پژوهش گران همواره در پی یافتن استانداردهای مناسب با ویژگی های منابع و بی در جهت سازماندهی هر چه بهتر آن ها بوده اند به همین جهت به ارزیابی و بررسی آن ها پرداخته اند.

4- روش شناسی پژوهش

روش این پژوهش تحلیل محتوا بود به صورتی که در ابتدا به شناسایی عناصر توصیف مرتبط با منابع اطلاعاتی وب در قالب کتاب شناختی یونی مارک و مارک 21 و سپس به تحلیل و مقایسه آن ها با یکدیگر پرداخته شد. برای انجام این پژوهش از قالب کتاب شناختی یونی مارک (International Federation of 2008 (Library Associations and Institutions (IFLA) و قالب کتاب شناختی مارک 21 (Library 2012 of Congress) استفاده شد. قابل ذکر است که شماری از عناصر دادهای موجود در قالب های کتاب شناختی یونی مارک و مارک 21 می توانند در پیشینه های منابع اطلاعاتی وب کاربرد داشته باشند، همان طور که در پیشینه های قالب های دیگر هم کاربرد دارند اما به دلیل تمرکز این پژوهش بر روی منابع اطلاعاتی، وب تنها عناصری که به طور خاص به این نوع از منابع اختصاص داشتند، مورد بررسی قرار گرفتند.

5- تجزیه و تحلیل یافته ها

برای پاسخ به پرسش نخست پژوهش به شناسایی فیلدها فیلدهای فرعی نشان گرها و نویسه های موجود در برچسب پیشینه مرتبط با منابع اطلاعاتی وب پرداخته شد برای پاسخ به این پرسش، علاوه دست نامه کامل هر دو قالب دست نامه آموزشی میلر برای فهرست نویسی منابع اینترنتی (Miller 2008)، راهنمای یونی مارک برای منابع الکترونیکی (International Federation of Library Associations and 2000). (Institutions (IFLA و جدول های تبدیل یونی مارک به مارک 21 (International Federation of Library Associations and Institutions (IFLA 2001 of Library Associations and Institutions) مرجع در نظر گرفته شد. در جدول 1 تنها فیلدها فیلدهای فرعی و نشان گرهایی که با محتوای منابع اطلاعاتی وب (به طور خاص) سازگار بودند، نشان داده شده است علاوه بر این معادل هر یک از عناصر در دو قالب در یک سطر قرار گرفت.

بررسی و مقایسه قابلیت های یونی مارک ... ۲۳۵

جدول ۱: جدول تطبیقی عناصر مربوط به منابع اطلاعاتی وب در قالب های یونی مارک و مارک ۲۱

یونی مارک				مارک ۲۱			
کد A شامل سیستمها و خدمات پیوسته هم می شود		نویسه ششم: نوع و کدرد/ کد A: فایبل رایانه ای		برجسب پیشینه		کد A شامل سیستمها و خدمات پیوسته هم می شود	
توضیحات	فیلدهای فرعی	نشانه ها		نام فیلد	شماره فیلد	توضیحات	فیلدهای فرعی
		اول	دوم				
عبارت "پیوسته" در فیلد فرعی Sb نشان می دهد که مدرک پیوند شده به شکل پیوسته می باشد.	Sb نام عام مواد	—	—	بلوک شناسه رابط	4—		
عبارت "پیوسته" در فیلد فرعی Sb نشان دهنده منابع وبی است	Sb نام عام مواد	—	—	عنوان و نام پدیدآور	200	عبارت "پیوسته" در فیلد فرعی Sh نشان دهنده منابع وبی است.	Sh: رسانه
	Sb نام عام مواد	—	—	عنوان قراردادی	500		Sh: رسانه
	Sb نام عام مواد	—	—	عنوان مشترک قراردادی	501		Sh: رسانه
عبارت "پیوسته" در فیلد فرعی Sn نشان دهنده منابع اطلاعاتی وب است	Sn اطلاعات مترقه	—	—	عنوان به منزله موضوع	605	عبارت "پیوسته" در فیلد فرعی Sg نشان دهنده منابع اطلاعاتی وب است	Sg: اطلاعات مترقه
006					کد 0 نشان دهنده شکل مدرک پیوسته است	نویسه 06: شکل مدرک	عناصر داده ای یا طول ثابت- ویژگی مواد اضافی
					فایبل نشان دهنده خدمات و سیستم های پیوسته است	نویسه 09: فایبل نشان دهنده نوع فایبل رایانه ای	
007	کد F نشان دهنده سیستم های پیوسته است.	Sa داده های کد شده برای منابع الکترونیکی نویسه 1: تعیین مواد خاص/ کد F	—	—	فیلد داده- های کد شده: منابع الکترونیکی	کد F نشان دهنده دسترسی از راه دور است	نویسه 01: تعیین مواد خاص
008	کد Z نشان دهنده خدمات و سرویس های پیوسته است	Sa داده های کد شده برای منابع الکترونیکی نویسه 0: نوع منبع الکترونیکی/ کد Z	—	—	فیلد داده- های کد شده: منابع الکترونیکی	کد Z نشان دهنده خدمات و سیستم های پیوسته است	نویسه 26: نوع فایبل رایانه ای
							نویسه 23: شکل مدرک

جدول 1 جدول تطبیقی عناصر مربوط به منابع اطلاعاتی وب در قالب های یونی مارک و مارک 21

۲۳۶ مدیریت منابع اطلاعاتی وب

256	مشخصات قابل رایانه- ای	—	—	—	متنقه خاص مواد مشخصات منابع الکترونیکی	230	اطلاعاتی همچون نوع قابل و اندازه آن در این فیلد فرعی وارد می‌شود	Sa بر حسب و اندازه قابل	اطلاعاتی همچون نوع قابل و اندازه آن در این فیلد فرعی وارد می‌شود
338	نوع حامل	—	—	—	—	—	بر اساس آردی‌ای عبارت "منابع" پیوسته "یکی از حامل‌های رایانه‌ای است.	Sa اصطلاح نوع حامل	اطلاعاتی همچون نوع قابل و اندازه آن در این فیلد فرعی وارد می‌شود
516	نوع قابل رایانه‌ای یا پادداشت داده	—	—	—	پادداشت نوع منابع الکترونیکی	336	عبارت "سیستم‌های پیوسته" در این فیلد فرعی نشان دهنده نوع منبع الکترونیکی است.	Sa نوع قابل رایانه‌ای یا پادداشت داده	عبارت "سیستم‌های پیوسته" در این فیلد فرعی نشان دهنده نوع منبع الکترونیکی است.
538	پادداشت جزئیات سیستم	—	—	—	پادداشت سیستم مورد نیاز	337	اطلاعات در مورد شکل دسترسی به منابع ویس در فیلد فرعی Sa آدرس دسترسی خودکار به مدرک در فیلد فرعی Su و بخشی از مواد شرح داده شده در S3 وارد می‌شود	Sa پادداشت جزئیات سیستم Si متن نمایش Su شناسگر متحدالشکل منبع S3 مواد معین	اطلاعات در مورد شکل دسترسی به منابع ویس در فیلد فرعی Sa آدرس دسترسی خودکار به مدرک در فیلد فرعی Su و بخشی از مواد شرح داده شده در S3 وارد می‌شود
753	جزئیات دسترسی سیستم به قابل رایانه- ای	—	—	—	جزئیات دسترسی فنی منابع الکترونیکی	626	—	Sa ساختار و الگوی دستگاهها Sb زبان‌های برنامه‌نویسی سیستم Sc عامل	—
856	دسترسی و تعیین محل الکترونیکی	روشن- های دسترسی	روشن- های دسترسی	تعریف نشده	دسترسی و تعیین محل الکترونیکی	856	Sa نام میزان؛ Sb شماره دسترسی؛ Sc فهرده‌سازی اطلاعات؛ Sd مسیر؛ Sf نام الکترونیکی؛ Sh پردازشگر درخواست؛ Si دستورالعمل؛ Sj بیت در هر تابه؛ Sk کلمه عبور؛ Sl ورود به سیستم؛ Sm برقراری تماس برای دریافت کمک؛ Sn نام محل میزان؛ So سیستم عامل؛ Sp درگاه؛ Sq نوع قالب الکترونیکی؛ Sr تنظیمات؛ Ss اندازه قابل؛ St تکرار پایانه‌ای؛ Su شناسگر متحدالشکل منبع؛ Sv ساعاتی که روش دسته‌ای قابل دستیابی است؛ Sw شماره کنترل پیشینه؛ Sx پادداشت غیر عمومی؛ Sy متن پیونده؛ Sz روش‌های دسترسی؛ S3 مواد مشخص شده؛ S6 پیونده؛ S8 پیونده فیلد و شماره توالی	Sa نام میزان؛ Sb شماره دسترسی؛ Sc فهرده‌سازی اطلاعات؛ Sd مسیر؛ Sf نام الکترونیکی؛ Sh پردازشگر درخواست؛ Si دستورالعمل؛ Sj بیت در هر تابه؛ Sk کلمه عبور؛ Sl ورود به سیستم؛ Sm برقراری تماس برای دریافت کمک؛ Sn نام محل میزان؛ So سیستم عامل؛ Sp درگاه؛ Sq نوع قالب الکترونیکی؛ Sr تنظیمات؛ Ss اندازه قابل؛ St تکرار پایانه‌ای؛ Su شناسگر متحدالشکل منبع؛ Sv ساعاتی که روش دسته‌ای قابل دستیابی است؛ Sw شماره کنترل پیشینه؛ Sx پادداشت غیر عمومی؛ Sy متن پیونده؛ Sz روش‌های دسترسی؛ S3 مواد مشخص شده؛ S6 پیونده؛ S8 پیونده فیلد و شماره توالی	روشن- های دسترسی

همان طور که در جدول 1 مشاهده می شود در هر دو قالب یکی از نویسه های برجسب پیشینه می تواند نشان دهنده منابع الکترونیکی باشد و در توضیحات ذیل کد آمده که شامل منابع پیوسته نیز هست. «برجسب پیشینه در ابتدای هر پیشینه قرار می گیرد و حاوی داده های مربوط به پردازش آن پیشینه است» (مارک ایران 1381 ، 28) و به طور غیر مستقیم برای کاربرد در تشخیص خود مدرک کتاب شناختی به کار می رود (International Federation of Library Associations and Institutions (IFLA). 2008 12).

همچنین نشان گرهای یکی از فیلدها (856) برای نمایش خصوصیات منابع وبی تعریف شده است. داده های جدول 1 نشان می دهد که در قالب مارک، 21 می توان در سیزده فیلد، و در قالب یونیمارک، می توان در دوازده فیلد اطلاعات مربوط به منابع وبی را وارد کرد.

برای پاسخ به پرسش دوم، پژوهش از خصوصیات تعریف شده برای منابع اطلاعاتی وب توسط هیری (1) (1996) استفاده و عناصر تعریف شده در دو قالب یونی مارک و مارک 21 مقایسه شد که نتایج آن در جدول 2 قابل مشاهده است.

عکس

همان‌طور که در جدول ۱ مشاهده می‌شود، در هر دو قالب، یکی از نویسه‌های برچسب پیشنهادی می‌تواند نشان دهنده منابع الکترونیکی باشد و در توضیحات ذیل کد آمده که شامل منابع پیوسته نیز هست. «برچسب پیشنهادی در ابتدای هر پیشنهاد قرار می‌گیرد و حاوی داده‌های مربوط به پردازش آن پیشنهاد است» (مارک ایران، ۱۳۸۱، ۲۸) و به‌طور غیرمستقیم برای کاربرد در تشخیص خود مدرک کتاب‌شناختی به‌کار می‌رود (International Federation of Library Associations and Institutions (IFLA). 2008. 12). همچنین، نشانگرهای یکی از فیلدها (۸۵۶) برای نمایش خصوصیات منابع وبی تعریف شده است. داده‌های جدول ۱ نشان می‌دهد که در قالب مارک ۲۱، می‌توان در سیزده فیلد، و در قالب یونی‌مارک، می‌توان در دوازده فیلد اطلاعات مربوط به منابع وبی را وارد کرد. برای پاسخ به پرسش دوم پژوهش، از خصوصیات تعریف شده برای منابع اطلاعاتی وب توسط هییری^۱ (۱۹۹۶) استفاده، و عناصر تعریف شده در دو قالب یونی‌مارک و مارک ۲۱ مقایسه شد که نتایج آن در جدول ۲ قابل مشاهده است.

جدول ۲: مقایسه خصوصیات منابع اطلاعاتی وب با عناصر موجود در قالب‌های یونی‌مارک و مارک ۲۱

ردیف	ویژگیها	مارک ۲۱	یونی‌مارک
۱	اطلاعات مربوط به مکان‌های مختلف یک منبع	فیلد ۸۵۶ و قابلیت تکرارپذیری آن برای ثبت مکان‌های مختلف یک منبع	فیلد ۸۵۶ و قابلیت تکرارپذیری آن برای ثبت مکان‌های مختلف یک منبع
۲	شیوه (های) دسترسی به منبع	نشانگر اول در فیلد ۸۵۶ و قابلیت تکرارپذیری برای بیش از یک روش دسترسی فیلد فرعی \$۲ در فیلد ۸۵۶ فیلد ۷۵۳	نشانگر اول در فیلد ۸۵۶ و قابلیت تکرارپذیری برای بیش از یک روش دسترسی فیلد فرعی \$۷ در فیلد ۸۵۶ فیلد ۶۲۶
۳	قالب‌های مختلف نسخه‌های مربوط به یک منبع (PDF, HTML, XML, ...)	فیلد فرعی \$q در فیلد ۸۵۶ فیلد فرعی \$y در فیلد ۸۵۶ فیلد ۳۳۸	فیلد فرعی \$q در فیلد ۸۵۶ فیلد فرعی \$b در بلوک ۴
۴	ثبت تغییرات ناشی از نبود ثبات	تکرارپذیری بودن فیلد ۸۵۶ در فیلدهای فرعی \$a, \$b, \$c	تکرارپذیری بودن فیلد ۸۵۶ در فیلدهای فرعی \$a, \$b, \$c
۵	روزآمد سازی منابع	_____	_____
۶	سطح توصیف	_____	_____
۷	اطلاعات مربوط به دسترسی و شرایط آن	فیلد ۵۳۸ فیلد ۲۵۶	فیلد ۳۳۷ فیلد ۲۳۰

داده‌های جدول ۲ نشان می‌دهد که برای ۵ مورد از خصوصیات بیان شده توسط هییری، حداقل یک عنصر در یونی‌مارک و مارک ۲۱ تعریف شده است. همچنین، برای نمایش دو مورد از خصوصیات نیز

1. Heery

جدول 2: مقایسه خصوصیات منابع اطلاعاتی وب با عناصر موجود در قالب های یونی مارک و مارک 21

داده های جدول 2 نشان می دهد که برای 5 مورد از خصوصیات بیان شده توسط هییری حداقل یک عنصر در یونی مارک 21 تعریف شده است. همچنین، برای نمایش دو مورد از خصوصیات نیز

برای پاسخ به پرسش سوم پژوهش با استفاده از داده‌های به دست آمده در جدول 1، به مقایسه عناصر تعریف شده در مارک 21 و یونیمارک پرداخته شد که نتایج این بررسی در جدول 3 نمایش داده شده است.

عکس

۲۳۸ مدیریت منابع اطلاعاتی وب

عنصری در مارک ۲۱ و یونی‌مارک یافت نشد.

برای پاسخ به پرسش سوم پژوهش، با استفاده از داده‌های به دست آمده در جدول ۱، به مقایسه عناصر تعریف شده در مارک ۲۱ و یونیمارک پرداخته شد که نتایج این بررسی در جدول ۳ نمایش داده شده است.

جدول ۳. مقایسه عناصر موجود در مارک ۲۱ و یونیمارک برای نمایش خصوصیات منابع اطلاعاتی وب

یونیمارک	مارک ۲۱
نمایش خصوصیات منابع اطلاعات وب در ۱۲ فیلد	نمایش خصوصیات منابع اطلاعات وب در ۱۳ فیلد
نمایش شکل مدرک در فیلدهای شناسه رابط	عدم نمایش شکل مدرک در فیلدهای شناسه رابط
عدم تعیین نوع مدرک در فیلدهای داده‌ای	تعیین نوع مدرک در فیلد داده‌ای (۰۰۸)
عدم تعیین نوع حامل بر اساس استاندارد آر.دی.ای.	تعیین نوع حامل بر اساس استاندارد آر.دی.ای.
در فیلد ۸۵۶ نشانگر دوم برای نمایش روابط تعریف شده است.	در فیلد ۸۵۶ نشانگر دوم برای نمایش روابط تعریف شده است.

داده‌های موجود در جدول ۳ نشان می‌دهد که قالب‌های یونی‌مارک و مارک ۲۱ برای نمایش خصوصیات منابع وبی در ۵ مقوله با یکدیگر متفاوت هستند. در قالب مارک ۲۱، در ۴ مقوله به نمایش خصوصیات و اختصاص عناصر جزئیتر در نمایش آنها توجه بیشتری در مقایسه با یونی‌مارک شده است. همین تفاوت در قالب یونی‌مارک در یک مقوله نسبت به مارک ۲۱ دیده می‌شود.

برای پاسخ به پرسش چهارم، به تجزیه و تحلیل نتایج به دست آمده در پرسش‌های اول و سوم پژوهش پرداخته شد که نتایج آن به قرار زیر است:

- فراوانی فیلدهای اختصاص داده شده به نمایش خصوصیات منابع وبی: با توجه به داده‌های موجود در جدول ۳، تفاوت مارک ۲۱ و یونی‌مارک در این مورد یک فیلد است. جدول ۱ نشان می‌دهد که فیلد فرعی \$b\$ در فیلدهای شناسه رابط یونی‌مارک نشان دهنده قالب مدرک پیوندی است. با توجه به ویژگی بلوک شناسه رابط در قالب یونی‌مارک، با نمایش قالب پیوسته در این فیلدها، علاوه بر مشخص شدن قالب پیوسته پیشینه‌های مرتبط با پیشینه اصلی، در هنگام بازیابی اطلاعات، قالب پیوسته مرتبط با مدرک بازیابی شده نیز مشخص خواهد شد.

در مارک ۲۱، فیلد ۰۰۶ مربوط به نمایش ویژگی مواد اضافی است. تفاوت این فیلد با فیلد ۰۰۸ این است که این فیلد برای مواردی که نمی‌توان آنها را در فیلد ۰۰۸ وارد کرد، آورده می‌شود. فیلد ۰۰۶ بیشتر برای مواردی که دارای چند نوع از یک گروه خصوصیت هستند استفاده می‌شود. به نظر می‌رسد فیلد ۰۰۶ جایگزینی برای فیلد فرعی \$b\$ در بلوک --۴ یونی‌مارک باشد، اگرچه فیلد فرعی \$b\$ علاوه بر نمایش قالب‌های دیگر مدرک مربوط، نمایش دهنده قالب تمام پیشینه‌های مرتبط (بر اساس انواع روابط موجود در میان آثار) است. علاوه بر این، در مارک ۲۱ بر اساس استاندارد آر.دی.ای، فیلد ۳۳۸ نشان دهنده نوع حامل است که در یونی‌مارک تعریف نشده است.

- تعیین نوع مدرک در فیلد داده‌ای: در مارک ۲۱، علاوه بر تعیین نوع فایل رایانه‌ای در فیلد داده‌ای،

مارک 21

داده های موجود در جدول 3 نشان می دهد که قالب های یونی مارک و مارک 21 برای نمایش خصوصیات منابع وبی در 5 مقوله با یکدیگر متفاوت هستند در قالب مارک ، 21، در 4 مقوله به نمایش خصوصیات و اختصاص عناصر جزئیتر در نمایش آن ها توجه بیشتری در مقایسه با یونی مارک شده است. همین تفاوت در قالب یونی مارک در یک مقوله نسبت به مارک 21 دیده می شود.

برای پاسخ به پرسش چهارم به تجزیه و تحلیل نتایج به دست آمده در پرسش های اول و سوم پژوهش پرداخته شد که نتایج آن به قرار زیر است:

- فراوانی فیلهای اختصاص داده شده به نمایش خصوصیات منابع وبی: با توجه به داده های موجود در جدول 3 تفاوت مارک 21 و یونی مارک در این مورد یک فیلد است جدول 1 نشان می دهد که فیلد فرعی b در فیلهای شناسه رابط یونی مارک نشان دهنده قالب مدرک پیوندی است. با توجه به ویژگی بلوک شناسه رابط در قالب یونی مارک با نمایش قالب پیوسته در این فیلهای، علاوه بر مشخص شدن قالب پیوسته پیشینه های مرتبط با پیشینه اصلی در هنگام بازیابی اطلاعات، قالب پیوسته مرتبط با مدرک بازیابی شده نیز مشخص خواهد شد.

در مارک ، 21 فیلد 006 مربوط به نمایش ویژگی مواد اضافی است تفاوت این فیلد با فیلد 008 این است که این فیلد برای مواردی که نمی توان آن ها را در فیلد 008 وارد کرد آورده می شود. فیلد 006 بیشتر برای مواردی که دارای چند نوع از یک گروه خصوصیت هستند استفاده می شود. به نظر می رسد فیلد 006 جایگزینی برای فیلد فرعی b در بلوک - یونی مارک ، باشد اگر چه فیلد فرعی b علاوه بر نمایش قالبهای دیگر مدرک مربوط نمایش دهنده قالب تمام پیشینه های مرتبط (بر اساس انواع روابط موجود در میان آثار) است. علاوه بر این در مارک 21 بر اساس استاندارد آر.دی.ای فیلد 338 نشان دهنده نوع حامل است که در یونی مارک تعریف نشده است.

- تعیین نوع مدرک در فیلد داده ای: در مارک 21، علاوه بر تعیین نوع فایل رایانه ای در فیلد داده ای

شکل مدرک نیز با کدی مشخص در نویسه 23 مشخص می شود باید توجه داشت که در مارک 21 مشخصات تمام انواع قالب ها در فیلد داده ای 008 وارد می شود، در صورتی که در یونی مارک بلوک --1 به اطلاعات کد شده اختصاص یافته و برای انواع مدارک دارای فیلد خاصی است. فیلد 135 مختص داده های کد شده برای منابع الکترونیکی است به همین دلیل نیازی به تعریف کدی برای شکل مدرک نبوده است.

- تعریف نشان گر برای نمایش: روابط در مارک ، 21 نشانگر دوم برای نمایش رابطه بین اطلاعات موجود در فیلد 856 و منبع توصیف شده در پیشینه است. این نشان گر ممکن است برای تولید نمایشی پیوسته یا نظم همگانی فیلدهای 856 استفاده شود (Network Development and MARC 2003 .Standards Office, Library of Congress) فیلد فرعی 13(1) نشان می دهد. به نظر می رسد نمایش روابط در فیلد 856 (حداقل برای منابع وبی) مناسب تر، باشد زیرا با توجه به این که تمام مشخصات منبع الکترونیکی (مانند مسیر شناس گر متحد الشكل منبع، و غیره) در فیلد 856 وارد می شود می توان قالب وبی هر منبع را (در صورت وجود) در پیشینه اصلی منبع وارد کرد و دیگر نیازی به تهیه پیشینه جدید برای قالب وبی آن نیست.

6- بحث و نتیجه گیری افزایش روزافزون منابع اطلاعاتی وب و آشفتگی موجود در بازیابی اطلاعات جامع و مانع، لزوم سازماندهی منابع اطلاعاتی وب را پر رنگ تر می سازد از آن جایی که منابع وبی سازماندهی شده جزئی از نظام های ذخیره و بازیابی کتابخانه ای خواهند بود لازم است برای سازماندهی این گونه منابع استانداردهای موجود در سازماندهی اطلاعات کتابخانه ای مورد بررسی قرار گیرند.

نتایج پژوهش حاضر نشان داد که قالب های مارک 21 و یونی مارک عناصری را به نمایش خصوصیات منابع وبی اختصاص داده اند این عناصر خصوصیات همچون، قالب نوع دسترسی، ملزومات دسترسی، و نشانی دسترسی را نمایش می دهند که با نتایج طاهری (1387) مبنی بر توانایی های مارک همخوانی داشت. علاوه بر این عناصر تعریف شده در این دو قالب توانایی نمایش سطح توصیف و روزآمدی را برای منابع اطلاعاتی وب. نداشتند نتایج این بخش از پژوهش با نتایج پژوهش وارد (2001) همخوانی نداشته و با نتایج پژوهش فارد و ریگو (2004) مبنی بر ناتوانی ابر داده ها در سازماندهی پیچیدگی های منابع الکترونیکی همخوانی دارد دو مورد بیان شده از مسائل اصلی در سازماندهی منابع وبی هستند زیرا همین خصوصیات منابع اطلاعاتی وب را از دیگر منابع اطلاعاتی متمایز و تصمیم گیری در چگونگی

ص: 239

1- نیز برای نمایش اطلاعات بیش تر در مورد وضعیتی است که در آن یک رابطه یک به یک وجود نداشته باشد. برای مثال نشان گر 1 نمایش گر وضعیتی است که مدرک توصیف شده در پیشینه کتاب شناختی الکترونیکی نیست اما یک نسخه الکترونیک از آن موجود و اطلاعات ثبت شده در فیلد 856 نشان دهنده نسخه الکترونیکی است (Library of Congress. 2003) همان طور که در بالا اشاره شد یونی مارک برای نمایش روابط از بلوک شناسه رابط استفاده می کند و قالب مدرک را در فیلد فرعی b

سازماندهی آن‌ها را چالش برانگیز نموده است. مقایسه دو قالب کتاب شناختی یونی مارک و مارک 21 نیز نشان داد که تفاوت چندانی در میان تعریف عناصر مربوط به قالب وبی وجود ندارد و تنها تفاوت‌هایی در شکل تعریف آن‌ها مشاهده شد به نظر می‌رسد در انتخاب یکی از دو قالب باید به وضعیت کنونی دو قالب در ایران میزان انطباق هر یک با استانداردها و الگوهای مطرح در سازماندهی اطلاعات وب، و همچنین وضعیت هر یک از آن‌ها در سازماندهی دیگر قالب‌های اطلاعاتی توجه شود.

از یافته‌های این پژوهش می‌توان چنین دریافت که قالب مارک می‌تواند به عنوان استاندارد برای ذخیره و بازیابی منابع اطلاعاتی وب به کار رود، با این حال مارک راه حلی برای چگونگی روزآمد نمودن و تعیین سطح توصیف منابع اطلاعاتی وب تهیه نکرده است پژوهشگران برای ویژگی سطح توصیف، استفاده از فیله‌های رابطه‌ای را پیشنهاد می‌کنند به این صورت که برای نمایش اجزای یک منبع وبی (مانند صفحه‌های یک وبگاه یا همان پیوندهای درونی) از فیله‌های رابطه‌ای کل و جزء و برای نمایش پیوندهای خارجی از فیله‌های رابطه‌ای که نشان دهنده روابطی همچون هم‌ارز هستند، استفاده شود. پژوهشگران برای ویژگی روزآمدی استفاده از فیله‌ی خاص برای ثبت تغییرات را پیشنهاد می‌دهند. این فیله شامل فیله‌های فرعی باشد که هر یک نشان دهنده یکی از تغییرات احتمالی (تغییر در آدرس، عنوان و غیره) در منابع اطلاعاتی وب و تعیین نشانگرها برای نمایش نوع تغییر (روزآمد شده، حذف شده، جایگزین شده و غیره) باشد با توجه به این که پژوهش حاضر به بررسی قالب مارک برای سازماندهی منابع اطلاعاتی وب به طور کلی پرداخته است انجام پژوهش‌هایی برای بررسی قالب مارک برای انواع منابع اطلاعاتی وب و ارزیابی توانایی‌های آن ضروری به نظر می‌رسد.

منابع

حاجی زین العابدینی محسن 1381 بررسی مسائل فهرست نویسی منابع اینترنتی و ارائه دست‌نامه پیشنهادی برای کتابخانه‌های ایران پایان‌نامه کارشناسی ارشد دانشگاه علوم پزشکی ایران. دانشکده مدیریت و اطلاع‌رسانی پزشکی

طاهری، سید مهدی 1387. مقایسه کارآیی طرح فراداده‌ای هسته‌دوئین و قالب فراداده‌ای مارک 21 در سازماندهی منابع اطلاعاتی شبکه جهانی وب فصلنامه کتابداری و اطلاع‌رسانی، 3(11): 139-158.

عمرانی ابراهیم 1390 چارچوب کتاب شناختی برای دوران دیجیتال خبرنامه انجمن کتابداری و اطلاع‌رسانی (27) 21-31

فتاحی، رحمت‌الله. 1380. چالش‌های سازماندهی دانش در قرن بیستم فصلنامه کتاب، 12(4): 59-83 کمیته ملی مارک ایران 1381 مارک ایران تهران کتابخانه ملی جمهوری اسلامی ایران

کوکبی، مرتضی، سمیه سادات آخشیک 1385 سازماندهی منابع اینترنتی: قواعد فهرست نویسی انگلو-امریکن یا عناصر فراداده‌های هسته دوئین؟ سازماندهی اطلاعات: رویکردها و راهکارهای نوین (مجموعه) مقالات اولین همایش انجمن کتابداری و اطلاع‌رسانی، ایران 16 و 17 اسفند 1385 (صص. 333-344) تهران: کتابدار

ملای، مقدم، گلناز محمد نعیم آبادی 1385 سازماندهی صفحات وب با استفاده از نظام های رده بندی کتابخانه ای سازماندهی اطلاعات: رویکرد ها و راهکار های نوین مجموعه مقالات اولین همایش انجمن کتابداری و اطلاع رسانی ایران، 16 و 17 اسفند 1385 (صص. 345-364). تهران: کتابدار.

Fard, S.E., A. Riggio. 2004. Medium or message? A new look at standards, structures, and schemata for managing electronic resources. *Library Hi Tech*, 22(2): 144-152

Gorman, M. 2004. Authority Control in the Context of Bibliographic Control in the Electronic Environment. *Cataloging Classification Quarterly*, 38(3-4). Retrieved on September 19, 2011, from <http://www.sba.unifi.it/ac/relazioni/gorman-eng.pdf>

Heery,R. 1996. Review of metadata formats. *Program: Electronic Library and Information Systems*, 30(4): 345-373

International Federation of Library Associations and Institutions (IFLA). 2000. UNIMARC guideline no.6: electronic resources. Retrieved on September 20, 2012, from [http:// archive.ifla.org/VI/3/p1996-1/guid6.htm](http://archive.ifla.org/VI/3/p1996-1/guid6.htm)

International Federation of Library Associations and Institutions (IFLA). 2001. UNIMARC to MARC 21 conversion specifications. Retrieved on September 12, 2012, from <http://www.loc.gov/marc/unimarc21.html>

International Federation of Library Associations and Institutions (IFLA). 2008. UNIMARC manual. München: K. G. Saur

Library of Congress 2003. MARC 21 format for bibliographic data: 856: Electronic location and access. Retrieved on September 19, 2012, from <http://www.loc.gov/marc/bibliographic/bd856.html>

Library of Congress 2012. MARC 21. Retrieved on September 19, 2012, from www.loc.gov/marc/translations.html

Miller, S. 2008. Rules and tools for cataloging internet resources (trainee manual). Retrieved on September 29, 2012, from <http://www.loc.gov/catworkshop/courses/cataloginginternet/pdf/ceig1-IM-FINAL.pdf>

Network Development and MARC Standards Office, Library of Congress 2003. MARC 21 formats: Guidelines for the use of field 856. Retrieved on September 29, 2012, from <http://www.loc.gov/marc/856guide.html>

- Ward, D. 2001. Internet resource cataloging: the SUNY Buffalo Libraries' response. OCLC Systems Services, 17(1): 19-26
- Weitz, J. 2006. Cataloging electronic resources: OCLC-MARC coding guidelines. Retrieved on September/ 20, 2012, from <http://www.oclc.org/support/documentation/worldcat/cataloging/electronicresources>

هدف پژوهش حاضر سنجش کیفیت محیط های رابط کاربر پایگاه های اطلاعاتی مجلات پیوسته تمام متن فارسی «سید»، «مگ ایران»، «نما متن»، «نورمگز»، «پژوهشگاه علوم و فناوری اطلاعات» و «مرکز منطقه ای» از طریق بررسی و مقایسه آن ها با معیارها و استانداردهای رعایت شده در سیاهه واریسی رابط کاربر پایگاه اطلاعاتی می باشد. روش پژوهش حاضر پیمایشی-توصیفی است گردآوری داده ها با استفاده از یک سیاهه واریسی محقق ساخته و از طریق مشاهده ی مستقیم انجام گرفته است. یافته ها نشان داد که پایگاه اطلاعاتی «پژوهشگاه علوم و فناوری اطلاعات» در مجموع با رعایت 35 مؤلفه (55/55 درصد) از مؤلفه های سیاهه واریسی برترین محیط رابط کاربری را از میان پایگاه های اطلاعاتی پیوسته مجلات تمام متن فارسی دارا می باشد محیط های رابط کاربری پایگاه های اطلاعاتی نور، مگز، مرکز منطقه ای و نامتن به ترتیب با میزان رعایت 34 مؤلفه (53/96 درصد)، 32 مؤلفه (50/79 درصد) و 28 مؤلفه (44/44 درصد) رتبه های دوم سوم و پنجم را به خود اختصاص دادند همچنین محیط کاربری دو پایگاه مگ ایران و سید به طور مساوی با میزان رعایت 31 مؤلفه (49/20 درصد) در جایگاه چهارم قرار گرفتند نتایج پژوهش نشان داد که محیط رابط کاربری پایگاه های تمام متن فارسی در امکانات نوینی که طی سال های اخیر پدیدار شده اند دچار کمبود می باشند از این رو پیشنهاد می شود طراحان پایگاه های مذکور به طراحی مجدد پایگاه های خود پرداخته و این امکانات و قابلیت ها را به پایگاه های خود بی افزایند.

کلیدواژه: رابط کاربر پایگاه های اطلاعاتی پیوسته مجلات تمام متن مجلات فارسی

صدیقه جعفر زاده (1) معصومه، پیروزفر، (2) عبدالحسین فرج پهلوی (3)

مقدمه

امروزه با رشد و گسترش شبکه جهانی اینترنت پایگاه های اطلاعاتی الکترونیک به عنوان مهم ترین و پر کاربرد ترین ابزار فراهم کننده ی دسترسی به اطلاعات تبدیل شده اند به طوری که پژوهشگران و محققان می توانند اطلاعات مورد نیاز خود را در میان انبوه اطلاعات از طریق این محمل های اطلاعاتی، به دقت و سرعت جستجو و بازیابی کنند نخستین نقطه برخورد کاربر با پایگاه های اطلاعاتی، محیط رابط کاربر است. در پایگاه های اطلاعاتی رابط کاربر عامل مهمی در تسهیل دستیابی کاربران به اطلاعات مورد نیاز خود محسوب می شود از این رو میزان ها و تولید کنندگان این پایگاه ها می کوشند تا ضمن در نظر گرفتن نیاز کاربران و با پیروی از اصول موجود عوامل و خصیصه های ضروری را جهت طراحی رابط کاربر پایگاه های اطلاعاتی خود شناسایی کرده و از این طریق دسترسی مؤثر کاربران نهایی را به اطلاعات موجود تضمین کنند (مهرداد و زاهدی 1386).

یکی از نمودهای عینی کیفیت استفاده کاربر از پایگاه های اطلاعاتی به ویژه پایگاه های اطلاعاتی پیوسته

ص: 243

1- دانشجوی کارشناسی ارشد کتابداری و اطلاع رسانی دانشگاه شهید چمران sedighehjafarzadeh88@gmail.com

2- دانشجوی کارشناسی ارشد کتابداری و اطلاع رسانی دانشگاه شهید چمران piruzfar@gmail.com

3- عضو هیئت علمی و مدیر گروه کتابداری و اطلاع رسانی دانشگاه شهید چمران farajpahlou@gmail.com

را می توان در قالبی با عنوان رابط کاربر مشاهده کرد (یمین فیروز 1382). در واقع رابط کاربر نخستین نقطه ی برخورد کاربر با پایگاه های اطلاعاتی است و به عنوان پلی ارتباطی بین انسان و سامانه ی اطلاعاتی عمل می کند. به همین دلیل مهم ترین هدف از طراحی رابط کاربر برآوردن رضایت کاربران و ایجاد تعامل بیش تر و بهتر بین کاربر و محیط های رایانه ای است (Shneiderman 1998 نقل در اعظمی و فتاحی 1388).

با توجه به این که رابط کاربر خوب باعث می شود کاربران مسیر خود را در پایگاه اطلاعاتی بهتر شناسایی کنند و تأثیر به سزایی در عملکرد آنان خواهد داشت (Bates 2002 ، نقل در علیجانی و دهقانی 1386)؛ بنابراین محیط رابط کاربر باید به کاربران کمک کند تا کلیدواژه ها و عبارت های مناسبی برای جستجوهای خود به کار ببرند منبع مورد نظر را از میان منابع اطلاعاتی انتخاب کنند، نتایج جستجو ها را درک و از چگونگی پیشرفت کارشان آگاه شوند تا بتوانند آسان تر و بهتر به اطلاعات مورد نیاز خود دست یابند (Hirst 1999)

در کشور ما نیز در دو دهه ی اخیر تعدادی پایگاه اطلاعاتی پیوسته ی تمام متن طراحی شده است که پژوهش گران و محققان بسیاری برای انجام کارهای علمی و پژوهشی و رفع نیازهای اطلاعاتی خود به آن ها مراجعه می کنند با توجه به نقش مهم رابط کاربر در دستیابی به محتوای پایگاه های اطلاعاتی و بالا بردن کیفیت خدمات ، آنان محیط رابط این پایگاه ها باید به گونه ای طراحی شده باشد که دسترسی مناسب و مؤثر به اطلاعات را برای کاربران تسهیل بخشد؛ به طوری که آنان بتوانند با درک روشنی از نحوه ی طراحی پایگاه اطلاعاتی ساختار اطلاعات و چگونگی ارائه اطلاعات به ، کاربران اطلاعات مورد نیاز خود را در پایگاه مورد نظر به سرعت جستجو و بازیابی کنند. بنابراین بررسی و مقایسه محیط رابط کاربر پایگاه های اطلاعاتی پیوسته مجلات تمام متن فارسی با معیارهای رعایت شده در پایگاه های اطلاعاتی معتبر پر کاربرد و استاندارد جهانی امری ضروری است این امر ضمن نشان دادن فاصله بین وضع موجود تا وضع مطلوب سبب بهبود کیفیت ارائه ی اطلاعات در پایگاه های اطلاعاتی مجلات فارسی می شود.

بیان مسأله و ضرورت پژوهش

پایگاه های اطلاعات علمی پیوسته مهم ترین و پرکاربرد ترین منابعی هستند که امروزه پژوهشگران جهت دسترسی به اطلاعات مورد نیاز خود به آن ها مراجعه می کنند در ایران نیز در همین راستا طراحی پایگاه های اطلاعاتی الکترونیکی آغاز شده و فعالیت های پراکنده ای توسط برخی سازمان ها در خصوص ارائه ی شکل الکترونیکی مجلات به صورت پیوسته یا ناپیوسته انجام گرفته است. از بین پایگاه هایی که به طور پیوسته شکل الکترونیکی نشریه های ایرانی را ارائه می دهند چهار پایگاه اطلاعاتی «سید» (1)، «مگ ایران»، «نما متن» و «نور مگز» جامعیت و پوشش موضوعی گسترده تری دارند. از طرفی، مدت زمان بیش تری از طراحی آن ها می گذرد بنابراین امکانات بالاتری نسبت به سایر پایگاه ها دارند و برای کاربران شناخته شده تر هستند (اسداللهی و نوکریزی 1389) در پژوهش حاضر علاوه بر این چهار پایگاه که صرفاً پایگاه های اطلاعاتی نشریات می باشند به بررسی محیط رابط کاربری پایگاه های اطلاعاتی مجلات

ص: 244

فارسی «پژوهشگاه علوم و فناوری اطلاعات» و «مرکز منطقه ای» (1) که دارای پایگاه مجلات و پایگاه های گوناگون دیگری نیز می باشند پرداخته می شود.

با توجه به استفاده گسترده کاربران از پایگاه های اطلاعاتی به عنوان منابع بازیابی اطلاعات و علی رغم تأکید متخصصان بر اهمیت رابط کاربر پایگاه ها و وب سایت ها پژوهش های انجام شده در این زمینه نشان می دهد که طراحان در پاره ای از موارد به دلیل در نظر نگرفتن شرایط و ویژگی های لازم از این مسئله غفلت کرده و در نتیجه مشکلاتی را در خصوص استفاده بهینه کاربران از خدمات خود به وجود آورده اند (Hansen 1998) همچنین بررسی ها نشان داد که به نظر می آید محیط رابط کاربری برخی پایگاه های اطلاعاتی پیوسته مجلات فارسی موجود در ایران در پاره ای موارد با استاندارد ها هماهنگی ندارند و دارای برخی کمبود ها می باشند از این رو وجود این مسئله پژوهش گران را به این امر واداشت تا بررسی و ارزیابی کنند که وضعیت محیط رابط کاربر پایگاه های اطلاعاتی پیوسته تمام متن فارسی در مقایسه با معیارهای رعایت شده در محیط رابط کاربری پایگاه های اطلاعاتی معتبر بین المللی چگونه است؟ و در این راستا نقاط قوت و ضعف این پایگاه ها کدامند؟ نتایج حاصل از این پژوهش می تواند سبب استفاده، بهتر افزایش موفقیت در جستجو ها و انجام بهتر و سریع تر بازیابی اطلاعات شود.

هدف پژوهش

هدف پژوهش حاضر سنجش کیفیت محیط های رابط کاربر پایگاه های اطلاعاتی مجلات پیوسته تمام متن فارسی «پژوهشگاه علوم و فناوری اطلاعات»، «سید»، «مرکز منطقه ای»، «مگ ایران»، «نما متن» و «نور مگز» از طریق بررسی و مقایسه آن ها با معیارهای رعایت شده در محیط رابط کاربر پایگاه های اطلاعاتی مجلات پیوسته تمام متن معتبر و برجسته بین المللی می باشد. به منظور دستیابی به این هدف نقاط قوت و ضعف موجود در محیط رابط کاربری این پایگاه ها مشخص می شود

با توجه به اهداف فوق سعی می شود به پرسش های زیر پاسخ داده شود:

1. وضعیت محیط رابط کاربری هر یک از پایگاه های اطلاعاتی مورد بررسی از لحاظ مقوله جستجو چگونه است؟
2. وضعیت محیط رابط کاربری هر یک از پایگاه های اطلاعاتی مورد، بررسی از لحاظ مقوله نمایش اطلاعات چگونه است؟
3. وضعیت محیط رابط کاربری هر یک از پایگاه های اطلاعاتی مورد بررسی از لحاظ مقوله خدمات چگونه است؟
4. وضعیت محیط رابط کاربری هر یک از پایگاه های اطلاعاتی مورد بررسی از لحاظ مقوله پیوندها چگونه است؟
5. وضعیت محیط رابط کاربری هر یک از پایگاه های اطلاعاتی مورد بررسی از لحاظ مقوله راهنمایی چگونه است؟

ص: 245

6. کدام یک از مقوله های کلی رابط کاربری پایگاه اطلاعاتی بیش ترین میزان رعایت را در پایگاه های اطلاعاتی مجلات الکترونیک پیوسته تمام متن فارسی دارد؟

7. کدام یک از پایگاه های اطلاعاتی مجلات الکترونیک پیوسته تمام متن فارسی بیش ترین انطباق را با معیارهای موجود در سیاهه واریسی محیط رابط کاربری پایگاه اطلاعاتی دارد؟

پیشینه ی پژوهش

سایین (2001) (1) رابط کاربر میزبان های پایگاه های اطلاعاتی را بر اساس پنج طبقه کلی عملکرد ارزیابی کرد و سپس خصوصیات بیش تری را که مربوط به هر طبقه بود به عنوان خصیصه های ضروری، مطلوب یا مورد نیاز مشخص کرد یافته های این مطالعه نشان داد که دایالوگ (2) با 32 امتیاز از مجموع 41 امتیاز بالاترین رتبه را دارد. او سی ال سی (3)، اوید (4) با 27 امتیاز در رتبه دوم پروکوئست و سیلور پلاتر (5) با 26 امتیاز در مقام سوم قرار دارند. نکسیس (6) با کسب 25 امتیاز در رتبه هفتم و ویلسون (7) با 20 امتیاز مقام آخر را کسب کرد. در پایان این پژوهش خصیصه های سیستم ایده آل مشخص شد بر این اساس دایالوگ به این سیستم بسیار نزدیک است اما هنوز فاقد تعدادی از خصیصه های مهم می باشد.

کراستینی (2004) (8) یک رابط کاربر گرافیکی جدید برای بازیابی مدارک به صورت سلسله مراتبی ارائه کرده و سپس به نحوی طراحی، اجرا و ارزیابی رابط کاربر پیشنهادی پرداخته است. نتایج این پژوهش حاکی از آن است که رابط کاربر پیشنهادی ابزارهای قدرتمند و مؤثری را برای جستجوی مدارک، رهیابی فهرست بازیابی شده و تصحیح جستجو فراهم می آورد.

شاید بتوان نخستین تلاش در ارزیابی رابط کاربر پایگاه های اطلاعاتی کتاب شناختی فارسی را پژوهش مهرداد و زاهدی (1386) دانست. آنان در پژوهش خود به بررسی و مقایسه محیط رابط کاربر دو میزبان داخلی (کتابخانه منطقه ای علوم و تکنولوژی (9) و پژوهشگاه اطلاعات و مدارک علمی ایران) با چهار میزبان خارجی (الزویر، امرالد، ابسکو و پروکوئست) ارائه دهنده ی پایگاه های اطلاعاتی به روش پیمایش تطبیقی پرداختند. با استفاده از سیاهه واریسی جامع با پنج خصیصه (خصیصه های کلی، جستجو، بازیابی، نمایش و کاربر پسندی) تلاش شد تا ضمن شناخت ویژگی ها و نقاط قوت و ضعف هر یک رابط کاربر میزبان های داخلی و خارجی با یکدیگر مقایسه شوند یافته های این پژوهش نشان داد در میزبان های داخلی به ترتیب کتابخانه منطقه ای علوم و تکنولوژی و پژوهشگاه اطلاعات و مدارک علمی ایران و در میزبان های خارجی به ترتیب ابسکو، پروکوئست امرالد و الزویر در پنج خصیصه مورد بررسی دارای

ص: 246

Sabin-1

Dialogweb-2

OCLC First Search-3

Ovid-4

Silver Platter-5

Nexis-6

WilsonWeb-7

9- منظور «مرکز منطقه ای اطلاع رسانی علوم و فناوری» می باشد

اعظمی و فتاحی (1388) به تعیین همخوانی محیط رابط پایگاه‌های اطلاعاتی اِبِسکو، امرالد پروکوئست و ساینس دایرکت با عناصر رفتار اطلاع‌یابی مدل «الیس» (1) که عبارتند از: شروع پیوندیابی، مرور تمایز، بازیابی و استخراج پرداختند نتایج نشان داد که عناصر «شروع» «پیوندیابی» و «تمایز» تا حدودی به وسیله‌ی برخی از محیط‌های رابط کاربر پایگاه‌های مورد بررسی حمایت می‌شوند، اما دیگر عناصر رفتار اطلاع‌یابی (تورق‌بازنگری و استخراج) در ساختار رابط کاربر این پایگاه‌ها لحاظ نشده‌اند. به طور کلی میزان مطابقت و همخوانی رابط کاربر پایگاه‌های اطلاعاتی با عناصر رفتار اطلاع‌یابی الیس در حد متوسط است بنابراین استفاده از این عناصر در طراحی و ارزیابی محیط رابط کاربری می‌تواند تأثیر زیادی بر بهینه‌شدن محیط رابط پایگاه‌های اطلاعاتی و در نتیجه بر فرآیند جستجو و بازیابی داشته باشد.

انتظاریان و فتاحی (1389) به تحلیل تبیین و شناسایی نقاط قوت و ضعف عناصر و ویژگی‌های مهم در محیط رابط پایگاه‌های اطلاعاتی مقاله‌های الکترونیکی مرکز منطقه‌ای اطلاع‌رسانی علوم و فناوری و پژوهشگاه اطلاعات و مدارک علمی پرداختند آن‌ها در این بررسی میزان همخوانی محیط رابط پایگاه‌های مورد بررسی با 10 مؤلفه نیلسن (2)، مشکلات اساسی محیط رابط این پایگاه‌ها و نیز تفاوت بین میزان درک کاربران متخصص و مبتدی را مورد سنجش قرار دادند. یافته‌ها نشان داد میزان همخوانی محیط رابط پایگاه‌های پژوهشگاه با 10 مؤلفه نیلسن به طور کلی در حد متوسط و در پایگاه‌های مرکز منطقه‌ای کمی بیش از حد متوسط است هر دو پایگاه در برخی از مؤلفه‌های مدل نیلسن دارای مشکلات اساسی هستند

اسداللهی و نوکاریزی (1389) طی پژوهشی به ارزیابی ساختار و محتوای پایگاه‌های اطلاعاتی الکترونیکی، سید مگ ایران و نما متن پرداختند. میزان مطابقت پایگاه سید، مگ ایران و نما متن با معیارهای ساختار به ترتیب 50/43 درصد، 54/70 درصد و 41/88 درصد برآورد شد. پایگاه نما متن بر خلاف دو پایگاه سید و مگ ایران به دلیل استفاده از زبان نمایه‌سازی کنترل شده، روند مشخص و منسجمی برای نمایه‌سازی مجله‌ها داشت. نتایج حاصل از بررسی ویژگی‌های منحصر به فرد ساختاری محتوایی پایگاه‌ها مبین آن بود که پایگاه سید از این نظر در رتبه‌ی نخست قرار دارد. به طور کلی هر سه پایگاه می‌توانند با توجه بیش‌تر به معیارهای مطرح شده در سیاهه واریسی بر قابلیت‌های ساختاری خود بیفزایند در بخش محتوا نیز افزایش جامعیت و توسعه‌ی پوشش، موضوعی استفاده از رویکرد ترکیبی یعنی زبان آزاد و کنترل شده، عدم گزینش در نمایه‌سازی افزایش پوشش گذشته‌نگر و کاهش دوره‌ی، روزآمد سازی در هر سه پایگاه توصیه می‌شود.

جمع‌بندی پیشنهادی

مرور پیشنهادها نشان می‌دهد که اهمیت بازیابی اطلاعات باعث گردیده پژوهش‌هایی بر روی پایگاه‌های

اطلاعاتی نشریات الکترونیک انجام گیرد. این پژوهش ها ابعاد مختلف این پایگاه ها از جمله رفتار اطلاع یابی در این پایگاه ها، گرافیک ساختار و محتوای پایگاه ها و به ویژه رابط کاربری این پایگاه ها را مورد بحث قرار داده اند اما نکته ای که حائز اهمیت است این است که تاکنون پژوهشی که رابط کاربری تمام پایگاه های شناخته شده مجلات فارسی را به صورت یک جا با معیارهای رعایت شده در چندین پایگاه بین المللی مورد ارزیابی قرار دهد انجام نشده است. از این رو در این پژوهش رابط کاربری پایگاه های مجلات فارسی با معیارهای پایگاه بین المللی مورد بررسی قرار می گیرد.

روش شناسی پژوهش

روش پژوهش حاضر پیمایشی- توصیفی است. جامعه ی مورد مطالعه در این پژوهش پایگاه های اطلاعاتی مجلات پیوسته تمام متن فارسی (پژوهشگاه علوم و فناوری اطلاعات)، «سید»، «مرکز منطقه ای»، «مگ ایران»، «نما متن» و «نور مگز» می باشد با توجه به اهداف پژوهش ابزار گردآوری داده ها با استفاده از یک سیاهه واریسی محقق ساخته است بدین ترتیب که محتوای سیاهه واریسی از معیارهای رعایت شده در محیط رابط کاربری پایگاه های اطلاعاتی مجلات پیوسته تمام متن معتبر و برجسته دنیا از قبیل «ساینس دایرکت» (1)، «پروکوئست» (2)، «امرالده» (3)، «ابسکو» (4) و «اشپرینگر» (5) و همچنین بررسی متون مرتبط استخراج و طراحی شد. سپس روایی آن توسط متخصصان این حوزه بررسی گردید آن گاه با مراجعه مستقیم به وب سایت پایگاه های اطلاعاتی مورد نظر بر اساس سیاهه واریسی به ارزیابی آن ها پرداخته شد. بر اساس معیارهای مطرح شده در سیاهه واریسی به ازای دارا بودن آن معیار عدد یک و در صورت نداشتن آن معیار عدد صفر منظور گردید برای تجزیه و تحلیل داده های گردآوری شده از طریق سیاهه واریسی نیز از آمار توصیفی (فراوانی درصد فراوانی و میانگین) استفاده شد.

یافته های پژوهش

در پاسخ به پرسش اول پژوهش باید گفت که پایگاه های اطلاعاتی سید و نور مگز به طور مساوی با میزان رعایت 11 مؤلفه (50 درصد) بیش ترین میزان رعایت معیارهای موجود در سیاهه واریسی را نسبت به پایگاه های اطلاعاتی دیگر داشته اند همچنین پایگاه های اطلاعاتی پژوهشگاه علوم و فناوری اطلاعات و مگ ایران و مرکز منطقه ای به طور مساوی 10 مؤلفه (45/45 درصد) و پایگاه اطلاعاتی نامتن 8 مؤلفه (36/36 درصد) از مؤلفه های مربوط به مقوله ی جستجو را در خود رعایت کرده اند. همچنین از مجموع 22 مؤلفه ی مربوط به مقوله ی جستجو مؤلفه های «جستجوی مجاورتی»، «ریشه سازی»، «توانایی اصلاح استراتژی جستجوی قبلی»، «محدود گره های مناسب جهت اصلاح جستجو قابلیت جستجوی متن آزاد»، «عملکردهای اصطلاحنامه» و «عملکرد مقالات مرتبط» در محیط رابط کاربری هیچ کدام از پایگاه های

ص: 248

ScienceDirect -1

ProQuest -2

Emerald -3

Ebsco -4

Springer -5

مورد بررسی رعایت نشده است. مؤلفه های «ذخیره جستجوها در حساب کاربری خود» و «تاریخچه جستجو» تنها در پایگاه اطلاعاتی پژوهشگاه علوم و فناوری اطلاعات رعایت شده است. به طور کلی می توان گفت محیط رابط کاربری پایگاه های اطلاعاتی مورد بررسی در این پژوهش به طور میانگین 45/45 درصد از مؤلفه های مربوط به جستجو را رعایت کرده اند (جدول 10)

عکس

سنجش رابط کاربری پایگاه های اطلاعاتی... ۲۴۹

مورد بررسی رعایت نشده است. مؤلفه های «ذخیره جستجوها در حساب کاربری خود» و «تاریخچه جستجو» تنها در پایگاه اطلاعاتی پژوهشگاه علوم و فناوری اطلاعات رعایت شده است. به طور کلی می توان گفت محیط رابط کاربری پایگاه های اطلاعاتی مورد بررسی در این پژوهش به طور میانگین ۴۵/۴۵ درصد از مؤلفه های مربوط به جستجو را رعایت کرده اند (جدول ۱۰).

جدول ۱. وضعیت رابط کاربری پایگاه های اطلاعاتی از لحاظ مقوله جستجو

مؤلفه های جستجو پایگاه های اطلاعاتی	مرکز منطقه ای	پژوهشگاه علوم و فناوری اطلاعات	نورمگز	مگ ایران	نماتن	سید	میانگین
جمع امتیازها	۱۰	۱۰	۱۱	۱۰	۸	۱۱	---
درصد	۴۵/۴۵	۴۵/۴۵	۵۰	۴۵/۴۵	۳۶/۳۶	۵۰	۴۵/۴۵

در پاسخ به پرسش دوم پژوهش می توان این طور بیان نمود که از لحاظ مقوله نمایش، محیط های رابط کاربری پایگاه پژوهشگاه علوم و فناوری اطلاعات با رعایت ۱۸ مؤلفه (۸۱/۸۱ درصد) بیشترین میزان رعایت و پایگاه مگ ایران با رعایت ۱۰ مؤلفه (۴۵/۴۵ درصد) کمترین میزان رعایت مؤلفه های مربوط به مقوله نمایش را داشته اند. همچنین از میان مؤلفه های مربوط به مقوله نمایش اطلاعات مؤلفه «نمایش آمار تعداد دانلود هر مقاله» در رابط کاربری هیچ کدام از پایگاه ها رعایت نشده است. به طور کلی یافته ها حاکی از آن است به طور میانگین پایگاه ها از لحاظ مقوله نمایش ۶۴/۳۸ درصد معیارها را رعایت کرده اند. نتایج در جدول ۲ قابل مشاهده است.

جدول ۲. وضعیت رابط کاربری پایگاه های اطلاعاتی از لحاظ مقوله نمایش اطلاعات

مؤلفه های نمایش پایگاه های اطلاعاتی	مرکز منطقه ای	پژوهشگاه علوم و فناوری اطلاعات	نورمگز	مگ ایران	نماتن	سید	میانگین
جمع امتیازها	۱۵	۱۸	۱۶	۱۰	۱۲	۱۴	---
درصد	۶۷/۱۸	۸۱/۸۱	۷۲/۷۲	۴۵/۴۵	۵۴/۵۴	۶۳/۶۳	۶۸/۳۸

پاسخ به پرسش سوم پژوهش: همان طور که از جدول ۳ مشخص است پایگاه های مگ ایران و نورمگز با رعایت ۵ مؤلفه (۶۲/۵۰ درصد)، در میزان هماهنگی با مؤلفه های مقوله خدمات در رتبه نخست قرار دارند. پایگاه نماتن و پژوهشگاه علوم و فناوری اطلاعات نیز با رعایت ۳ مؤلفه (۳۷/۵۰ درصد)، پایگاه مرکز منطقه ای با رعایت ۲ مؤلفه (۲۵ درصد)، پایگاه سید به ترتیب در رتبه های دوم، سوم و

جدول ۱. وضعیت رابط کاربری پایگاه های اطلاعاتی از لحاظ مقوله جستجو

در پاسخ به پرسش دوم پژوهش می توان این طور بیان نمود که از لحاظ مقوله نمایش محیط های رابط کاربری پایگاه پژوهشگاه علوم و فناوری اطلاعات با رعایت 18 مؤلفه (81/81 درصد) بیش ترین میزان رعایت و پایگاه مگ ایران با رعایت 10 مؤلفه (45/45) درصد کم ترین میزان رعایت مؤلفه های مربوط به مقوله نمایش را داشته اند همچنین از میان مؤلفه های مربوط به مقوله نمایش اطلاعات مؤلفه ی «نمایش آمار تعداد دانشجو هر مقاله در رابط کاربری هیچ کدام از پایگاه ها رعایت نشده است. به طور کلی یافته ها حاکی از آن است به طور میانگین پایگاه ها از لحاظ مقوله نمایش 64/38 درصد معیارها را رعایت کرده اند. نتایج در جدول 2 قابل مشاهده است

جدول 2. وضعیت رابط کاربری پایگاه های اطلاعاتی از لحاظ مقوله نمایش اطلاعات

پاسخ به پرسش سوم پژوهش همان طور که از جدول 3 مشخص است پایگاههای مگ ایران و نورمگز با رعایت 5 مؤلفه (62/50 درصد) در میزان هماهنگی با مؤلفه های مقوله خدمات در رتبه نخست قرار دارند پایگاه نما متن و پژوهشگاه علوم و فناوری اطلاعات نیز با رعایت 3 مؤلفه (37/50 درصد)، پایگاه مرکز منطقه ای با رعایت 2 مؤلفه (25 درصد) پایگاه سید به ترتیب در رتبه های دوم، سوم و

ص: 249

چهارم قرار گرفتند همچنین از میان مؤلفه های مقوله خدمات مؤلفه های «امکان ساخت پرونده شخصی (شخصی سازی محتوی)» و سهولت خروج از سیستم به طور مساوی با فراوانی 4 و درصد فراوانی 80 بیشترین میزان رعایت را داشته اند مؤلفه های «استخراج اطلاعات کتابشناختی برای نرم افزارهای مدیریت ارجاعات» و «ارائه خدمات آگهی رسانی جاری (آلرت نشریات، آلرت جستجو)» در محیط های رابط کاربری هیچ کدام از پایگاه های اطلاعاتی مورد بررسی رعایت نشده است. میانگین امتیاز پایگاه ها در این مقوله نیز 37/49 درصد به دست آمد (جدول 3)

عکس

۲۵۰ مدیریت منابع اطلاعاتی وب

چهارم قرار گرفتند. همچنین از میان مؤلفه های مقوله خدمات، مؤلفه های «امکان ساخت پرونده شخصی (شخصی سازی محتوی)» و «سهولت خروج از سیستم» به طور مساوی با فراوانی ۴ و درصد فراوانی ۸۰ بیشترین میزان رعایت را داشته اند. مؤلفه های «استخراج اطلاعات کتابشناختی برای نرم افزارهای مدیریت ارجاعات» و «ارائه خدمات آگهی رسانی جاری (آلرت نشریات، آلرت جستجو)» در محیط های رابط کاربری هیچ کدام از پایگاه های اطلاعاتی مورد بررسی رعایت نشده است. میانگین امتیاز پایگاه ها در این مقوله نیز ۳۷/۴۹ درصد به دست آمد (جدول ۳).

جدول ۳. وضعیت رابط کاربری پایگاه های اطلاعاتی از لحاظ مقوله خدمات

مؤلفه های خدمات	میانگین	سید	نماتن	مگ ایران	نورمگز	پژوهشگاه علوم و فناوری اطلاعات	مرکز منطقه ای	پایگاه های اطلاعاتی
جمع امتیازها	---	۰	۳	۵	۵	۳	۲	جمع امتیازها
درصد	۳۷/۴۹	۰	۳۷/۵۰	۶۲/۵۰	۶۲/۵۰	۳۷/۵۰	۲۵	درصد

در پاسخ به سؤال چهارم می توان گفت پایگاه های مرکز منطقه ای، مگ ایران، نماتن و سید هر یک با رعایت ۳ مقوله (۶۰ درصد) مشترکاً در رتبه نخست و پایگاه های پژوهشگاه علوم و فناوری اطلاعات و نورمگز به طور مساوی با رعایت ۲ مقوله (۴۰ درصد) در رتبه دوم قرار دارند. پایگاه ها از لحاظ مقوله پیوندها به طور میانگین ۵۲/۳۳ درصد از مؤلفه ها رادر محیط رابط کاربری خود رعایت کرده اند. نتایج در جدول ۴ آمده است.

جدول ۴. وضعیت رابط کاربری پایگاه های اطلاعاتی از لحاظ مقوله پیوندها

مؤلفه های پیوندها	میانگین	سید	نماتن	مگ ایران	نورمگز	پژوهشگاه علوم و فناوری اطلاعات	مرکز منطقه ای	پایگاه های اطلاعاتی
جمع امتیازها	---	۳	۳	۳	۲	۲	۳	جمع امتیازها
درصد	۳۷/۴۹	۶۰	۶۰	۶۰	۴۰	۴۰	۶۰	درصد

پاسخ سؤال پنجم: همان گونه که از جدول ۵ مشاهده می شود پایگاه مرکز منطقه ای با رعایت ۴ مقوله (۶۶/۶۶ درصد) در رتبه نخست، پایگاه های نورمگز، مگ ایران و سید با رعایت ۳ مقوله مشترکاً در رتبه دوم، پایگاه پژوهشگاه علوم و فناوری اطلاعات با رعایت ۲ مقوله (۳۳/۳۳ درصد) در رتبه سوم و پایگاه نماتن تنها با رعایت ۱ مقوله (۱۶/۶۶ درصد) در رتبه چهارم قرار دارد. همچنین از لحاظ مقوله راهنمایی مؤلفه های «پیام های اخطاری قابل درک» و «نقشه سایت» به ترتیب با درصدهای فراوانی ۱۰۰ و

جدول 3. وضعیت رابط کاربری پایگاه های اطلاعاتی از لحاظ مقوله خدمات

در پاسخ به سؤال چهارم می توان گفت پایگاه های مرکز منطقه ای مگ، ایران، نامتن و سید هر یک با رعایت مقوله (60 درصد) مشترکاً در رتبه نخست و پایگاه های پژوهشگاه علوم و فناوری اطلاعات و نورمگز به طور مساوی با رعایت 2 مقوله (40 درصد) در رتبه دوم قرار دارند پایگاه ها از لحاظ مقوله پیوندها به طور میانگین 53/33 درصد از مؤلفه ها را در محیط رابط کاربری خود رعایت کرده اند. نتایج در جدول 4 آمده است.

جدول 4. وضعیت رابط کاربری پایگاه های اطلاعاتی از لحاظ مقوله پیوندها

پاسخ سؤال پنجم همان گونه که از جدول 5 مشاهده می شود پایگاه مرکز منطقه ای با رعایت 4 مقوله (6666) درصد در رتبه نخست پایگاه های نورمگز مگیران و سید با رعایت 3 مقوله مشترکاً در رتبه دوم پایگاه پژوهشگاه علوم و فناوری اطلاعات با رعایت 2 مقوله (33/33 درصد) در رتبه سوم و پایگاه نامتن تنها با رعایت 1 مقوله (1666) (درصد در رتبه چهارم قرار دارد هم چنین از لحاظ مقوله راهنمایی مؤلفه های پیام های اختطاری قابل درک و نقشه سایت به ترتیب با درصدهای فراوانی 100 و

ص: 250

بیشترین و کمترین میزان رعایت را در رابط کاربری پایگاه‌های مورد بررسی داشته‌اند. به طور میانگین 44/44 درصد مؤلفه‌های مربوط به مقوله راهنمایی توسط پایگاه‌های مورد بررسی رعایت شده‌اند.

عکس

سنجش رابط کاربری پایگاه‌های اطلاعاتی... ۲۵۱

• بیشترین و کمترین میزان رعایت را در رابط کاربری پایگاه‌های مورد بررسی داشته‌اند. به طور میانگین ۴۴/۴۴ درصد مؤلفه‌های مربوط به مقوله راهنمایی توسط پایگاه‌های مورد بررسی رعایت شده‌اند.

جدول ۵. وضعیت رابط کاربری پایگاه‌های اطلاعاتی از لحاظ مقوله راهنمایی

مؤلفه‌های راهنمایی پایگاه‌های اطلاعاتی	مرکز منطقه‌ای	پژوهشگاه علوم و فناوری اطلاعات	نورمگز	مگ‌ایران	نماتن	سید	میانگین
جمع امتیازها	۴	۲	۳	۳	۱	۳	---
درصد	۶۶/۶۶	۳۳/۳۳	۵۰	۵۰	۱۶/۶۶	۵۰	۴۴/۴۴

پاسخ به پرسش ششم پژوهش: همان‌گونه که از جداول ۱ تا ۵ ملاحظه می‌گردد بیشترین امتیاز و درصد فراوانی میزان رعایت مقوله‌های مورد بررسی در محیط رابط کاربری پایگاه‌های اطلاعاتی مرکز منطقه‌ای، پژوهشگاه علوم و فناوری اطلاعات، نورمگز و سید به ترتیب با امتیازهای ۱۵، ۱۶، ۱۸، ۱۴ و درصدهای ۶۷/۱۸، ۸۱/۸۱، ۷۲/۷۲ و ۶۳/۶۳ مربوط به مقوله‌ی نمایش و می‌باشد. همچنین بیشترین امتیاز و درصد فراوانی میزان رعایت مقوله‌های مورد بررسی در محیط رابط کاربری پایگاه‌های اطلاعاتی مگ‌ایران و نماتن به‌طور مساوی با ۳ امتیاز و ۶۰ درصد مربوط به مقوله پیوندها است.

پاسخ به پرسش هفتم پژوهش: همان‌طور که در جدول ۶ مشاهده می‌شود محیط رابط کاربری پایگاه اطلاعاتی پژوهشگاه علوم و فناوری اطلاعات در مجموع با میزان رعایت ۳۵ مؤلفه (۵۵/۵۵ درصد) بیشترین میزان رعایت مقوله‌های موجود در سیاهه‌وارسی را داشته است و رتبه نخست را به خود اختصاص داده است. سپس محیط‌های رابط کاربری پایگاه‌های نورمگز، مرکز منطقه‌ای، سید و مگ‌ایران به‌طور مساوی و نماتن در رتبه‌های دوم تا چهارم قرار داشتند.

جدول ۶. وضعیت کلی رابط کاربری پایگاه‌های اطلاعاتی

مؤلفه کلی پایگاه‌های اطلاعاتی	۱	۲	۳	۴	۵	۶	۷
مرکز منطقه‌ای	۸	۱۵	۲	۳	۴	۳۲	۵۰/۷۹
پژوهشگاه علوم و فناوری اطلاعات	۱۰	۱۶	۲	۲	۳	۳۵	۵۵/۵۵
نورمگز	۱۱	۱۶	۵	۲	۳	۳۴	۵۳/۹۶
مگ‌ایران	۱۰	۱۰	۵	۳	۳	۳۱	۴۹/۲۰
نماتن	۸	۱۲	۳	۳	۱	۲۷	۴۲/۸۵
سید	۱۱	۱۴	۰	۳	۳	۳۱	۴۹/۲۰

جدول ۵. وضعیت رابط کاربری پایگاه‌های اطلاعاتی از لحاظ مقوله راهنمایی

پاسخ به پرسش ششم پژوهش: همان گونه که از جداول 1 تا 5 ملاحظه می گردد بیش ترین امتیاز و درصد فراوانی میزان رعایت مقوله های مورد بررسی در محیط رابط کاربری پایگاه های اطلاعاتی مرکز منطقه ای، پژوهشگاه علوم و فناوری اطلاعات نورمگز و سید به ترتیب با امتیازهای 15، 18، 16، 14 درصدهای 18/18، 81/81، 72/72 و 63/63 مربوط به مقوله ی نمایش و می باشد. همچنین بیشترین امتیاز و درصد فراوانی میزان رعایت مقوله های مورد بررسی در محیط رابط کاربری پایگاه های اطلاعاتی مگ ایران و نما متن به طور مساوی با 3 امتیاز و 60 درصد مربوط به مقوله پیوندها است.

پاسخ به پرسش هفتم پژوهش: همان طور که در جدول 6 مشاهده می شود محیط رابط کاربری پایگاه اطلاعاتی پژوهشگاه علوم و فناوری اطلاعات در مجموع با میزان رعایت 35 مؤلفه (55/55 درصد) بیش ترین میزان رعایت مقوله های موجود در سیاهه واری را داشته است و رتبه نخست را به خود اختصاص داده است. سپس محیط های رابط کاربری پایگاههای نورمگز مرکز منطقه ای، سید و مگایران به طور مساوی و نامتن در رتبه های دوم تا چهارم قرار داشتند.

جدول 6. وضعیت کلی رابط کاربری پایگاه های اطلاعاتی

ص: 251

این پژوهش با هدف سنجش رابط کاربری پایگاه های اطلاعاتی پیوسته مجلات تمام متن فارسی انجام شد. نتایج نشان می دهد که از لحاظ مقوله ی جستجو پایگاه ها در سطح مشابه ای می باشند و در حد متوسط (میانگین 45/45 درصد) معیارها را رعایت کرده اند. از لحاظ مقوله ی نمایش نیز پایگاه ها در وضعیت خوبی میانگین (64/38) قرار دارند از این لحاظ رابط کاربری پایگاه پژوهشگاه علوم و فناوری اطلاعات با میزان رعایت 81/81 درصد و پایگاه مگ ایران با میزان رعایت 45/45 درصد به ترتیب بیش ترین و کم ترین میزان هماهنگی را با مقوله های موجود در سیاهه واری داشته اند. رابط کاربری پایگاه ها در میزان رعایت مقوله خدمات کم ترین هماهنگی را نسبت به سایر مقوله ها داشته اند؛ به طوری که میانگین رعایت معیارهای این مقوله 37/49 درصد و در حد ضعیف محاسبه شده است. در ارزیابی میزان رعایت مؤلفه های مقوله پیوندها نیز رابط کاربری پایگاه ها در سطح متوسطی قرار داشتند و به طور میانگین 53/33 درصد معیارها را رعایت نموده اند از لحاظ میزان رعایت مؤلفه های مقوله راهنمایی به جز پایگاه نما متن که تنها 16/66 درصد معیارهای این مقوله را رعایت نموده است؛ رابط کاربری سایر پایگاه های مورد بررسی در وضعیت مشابه و متوسطی (میانگین 44/44 درصد) قرار داشتند.

همچنین نتایج دیگری که از این پژوهش حاصل شد این است که پایگاه اطلاعاتی «پژوهشگاه علوم و فناوری اطلاعات» با بیشترین میزان رعایت، مؤلفه ها برترین محیط رابط کاربری را از بین محیط رابط کاربری پایگاه های اطلاعاتی پیوسته مجلات تمام متن فارسی دارا می باشد در نهایت می توان اظهار نمود که محیط رابط کاربری پایگاههای مورد بررسی از لحاظ میزان رعایت معیارهای موجود در سیاهه واری در حد «متوسطی» می باشند و این نتیجه یافته های پژوهش های انتظاریان و فتاحی (1389)، اعظمی و فتاحی (1388) و اسداللهی و نوکریزی (1389) را تأیید می نماید به طور کلی بهتر است هر شش پایگاه نسبت به گنجاندن آن دسته از معیارهای مطرح شده در سیاهه واری که فاقد آن هستند، اقدام کنند.

در نهایت پژوهشگران بر اساس نتایجی که از ارزیابی رابط کاربری پایگاه ها به دست آمد به طراحان محیط رابط کاربری پایگاه های اطلاعاتی مجلات تمام متن فارسی پیشنهاد می دهند که جهت استفاده بهتر، افزایش موفقیت در جستجوها و انجام بهتر و سریعتر بازیابی اطلاعات از طریق محیط رابط کاربری و در نتیجه جذب کاربران بیشتر برخی نکات را که در محیط رابط کاربری آن ها به ندرت استفاده شده یا اصلاً استفاده نشده است رعایت کنند. این نکات عبارتند از: در مقوله جستجو به طراحی آیتم های «محدودگرهای مناسب»، «عملکرد اصلاح نامه»، «عملکرد مقالات مرتبط»، «عملکرد پیشنهاد کلیدواژه های مرتبط»، «ذخیره جستجوها در حساب کاربری» و «تاریخچه جستجو»؛ در زمینه نمایش اطلاعات به طراحی آیتم های «نشانه دار کردن نتایج» و «ارسال نتایج به ایمیل»؛ از لحاظ مقوله خدمات پایگاه ها «امکان استخراج اطلاعات کتابشناختی به نرم افزارهای مدیریت ارجاعات» «نشان دادن مقالات استناد کننده» و «خدمات آگاهی رسانی جاری»؛ در زمینه پیوندها نیز «امکان فرایبوند به عناوین مشابه» و «فرایبوند به دیگر مقالات یک نویسنده» پردازند. بنابراین پیشنهاد می شود که طراحان پایگاه ها به طراحی مجدد پایگاه های خود پرداخته و به ویژه خدمات جدید معرفی شده در این پژوهش را در محیط رابط کاربری پایگاه های اطلاعاتی خود اعمال نمایند.

پیشنهاد می شود پایگاه های مذکور را از نظر جنبه های خاص شامل «دامنه پوشش موضوعی و زمانی مجلات»، «نحوه سازمان دهی اطلاعات» و غیره مورد پژوهش قرار دهند همچنین مطالعاتی نیز بر پایگاه های تهیه شده در کشور جز پایگاه های مجلات از قبیل پایگاه های اطلاعات کتابشناختی پایگاه های مقالات کنفرانس ها انجام شود.

منابع

اسد اللهی، زهرا، محسن نو کاریزی. 1389. ارزیابی ساختار و محتوای پایگاه های اطلاعاتی الکترونیکی نشریات ایرانی کتابداری و اطلاع رسانی. 13(2)

در [http://www.aqlibrary.ir/index.php?module=TWArticlesfile=indexfunc=view&pubarticlesdid=882pid=10.\(91/5/3](http://www.aqlibrary.ir/index.php?module=TWArticlesfile=indexfunc=view&pubarticlesdid=882pid=10.(91/5/3) (دسترسی در

اعظمی، محمد رحمت الله فتاحی 1388. تطابق رابط گرافیکی کاربر پایگاه های اطلاعاتی با مدل رفتار ابی ایس علوم و فناوری اطلاعات. 25 (2): 264-247.

انتظاریان، ناهید، رحمت الله فتاحی. 1389. مبانی طراحی رابط کاربر مبتنی بر شناخت ویژگی ها، ادراک و رفتار کاربران. کتابداری و اطلاع رسانی. 13 (2).

در [http://www.aqlibrary.ir/index.php?module=TWArticlesfile=indexfunc=view&pubarticlesdid=878pid=10.\(91/4/24](http://www.aqlibrary.ir/index.php?module=TWArticlesfile=indexfunc=view&pubarticlesdid=878pid=10.(91/4/24) (دسترسی در

علیجانی، رحیم، لیلا دهقانی 1386 مقایسه ی رابط کاربر پایگاه های اطلاعاتی کتاب مدار بین المللی فصلنامه کتاب 72: 252-233
مهرداد، جعفر، لیلا- دهقانی 1385 معیارهای ارزیابی رابط های کاربر نسخه های مختلف پایگاه های اطلاعاتی مجله کتابداری بهار و تابستان 77-95.

مهرداد جعفر زهره زاهدی. 1386 بررسی و مقایسه رابط کاربر دو میزبان داخلی کتابخانه منطقه ای علوم و تکنولوژی و پژوهشگاه اطلاعات و مدارک علمی ایران با چهار میزبان خارجی Emerald، Ebsco، Proquest Elsevier کتابداری و اطلاع رسانی 10 (3): 124-107.

یمین فیروز، موسی 1382 ویژگی ها و عناصر تشکیل دهنده رابط کاربر در وب سایت ها فصلنامه کتاب. 14: 168-159

Craestini, F. 2004. A Graphical user interface for the retrieval of hierarchically structured documents. Information Processing Management. 40 (2): 269-289

Hansen, P. 1998. Evaluation of IR user interface implications for user interface design. www.)
Hb.Se/bhs/ith/2-98/ph.htm (Available at 12/6/2012

Hirst, S.J. 1999. HyperLib Deliverable 2.1.1: The Use of Icons in a Multilingual OPAC Interface. Hyperlib
.Electronic Document Store, University of Antwerp–University of Loughborough

<http://lib.ua.ac.be/MAN/WP211/root.htm>(Available at 8/5/2012)

Sabin–Kidiss, L. 2001. Assessing the functionality of web–based versions of traditional search engines
..Online. 25 (2): 18–24

ص: 253

طبق قانون حق مؤلف کشور فرانسه مصوب اول اوت 2006، کتابخانه ملی فرانسه (BNF) از این پس کتابخانه (وظیفه گردآوری و حفاظت اینترنت فرانسه را بر عهده دارد. این کتابخانه «مدل تلفیقی» را برای آرشیو وب ایجاد کرده است. این مدل، متشکل از خزش های فراگیر دامنه fr و خزش های کانونی و اسپاری های الکترونیکی است.

کتابخانه ملی فرانسه، به برکت همکاری پژوهشی با آرشیو اینترنت از سال 2004 هر سال به اجرای چهار خزش فراگیر مبادرت کرده است. آخرین آن ها، با ویژگی های متفاوت چشمگیری ایجاد شده است و از مهم ترین این ویژگی ها کاربرد فهرست همه جانبه ای از اسامی دامنه های fr، بود که توسط آفینک انجمن همکاری نام گذاری اینترنت، فرانسه برای ثبت (fr) بعد از امضای توافق نامه ای میان دو سازمان در سپتامبر 2007 به کتابخانه ملی فرانسه ارائه شد.

گزینش های ماهرانه قبل و حین، خزش نقش تعیین کننده های در شکل آینده مجموعه خواهند داشت، بنابراین تصمیم هایی که باید بر طبق ساختار قانونی و معنوی در مدت انجام خزش اتخاذ شوند عبارت اند از: برای بیان اف این مجموعه شامل 5 قرن سنت گذشته و اسپاری قانونی است. برای تعیین پیامدها و نتایج راه حل های فنی موجود، قصد داریم نتایج آخرین خزش این کتابخانه را تحلیل و با برداشت های سال های گذشته مقایسه کنیم به علاوه این مطالعات سودمندی تلاش های ما را برای توصیف وب 2007 فرانسه تأیید می کند.

کلید واژه ها: آرشیو سازی، وب قانون و اسپاری اینترنت کتابخانه ملی فرانسه بیان اف، کنسرسیوم بین المللی حفاظت، اینترنت IIPC آرشیو اینترنت میراث دیجیتال

راهبردهایی برای گردآوری دامنه ملی (1)

نوشته: فرانس لاس فارگوس کلمنت کیوری، برت وندلاند (2)

ترجمه: سودابه نوذری (3)

1. قلمرو فرانسه

1,1 تعریف حوزه قانون و اسپاری

در اول اوت 2006، قانون حق مؤلف جدید در مجلس فرانسه به رأی گذاشته شد. مفاد این قانون، که کتابخانه مدت ها انتظارش را می کشید، قانون واسپاری را به اینترنت کشاند. قانون واسپاری، به هر ناشر تکلیف می کند نسخه هایی از تولیداتش را به کتابخانه ملی ارسال کند این، قانون که نخستین بار در 1537 برای منابع چاپی تصویب شده بود با گذشت قرن ها به اشکال متفاوت تولیدات فرهنگی انتشاراتی جدید از منقوشات گرفته تا نرم افزارها و بازی های ویدئویی تعلق گرفت به موازات گسترش شبکه جهانی وب، به عنوان مکانی مطلوب برای ایجاد دانش و اطلاعات لازم بود برای نهادهای حافظ میراث ملی فرانسه، چارچوبی قانونی برای سازماندهی منابع مورد حفاظت شان ارائه شود.

قانون فوق درباره قلمرو اینترنت فرانسه شفاف نیست و برای روشن شدن این مسئله در آینده ای

ص: 255

Legal deposit of the French Web: harvesting strategies for a national domain –1

France Lasfargues, Clément Oury, and Bert Wendland –2

3- عضو هیئت علمی سازمان اسناد و کتابخانه ملی ایران

نزدیک صدور فرمانی انتظار می رود در این اثنا، کتابخانه خط مشی خود را درباره دامنه ملی برخط به پیش می برد که احتمالاً با، فرمان زمانی که تأیید شود سازگار باشد، کتابخانه به منظور پالایش حوزه سیاست آرشیوسازی وب خود پنج قرن و اسپاری قانونی و تجربه جدیدتر چالش های فنی گردآوری وب را با یکدیگر تلفیق کرده است بنابراین روش ما هم عملی است (یعنی موافق با ابزار اکتشافی برداشت حجمی) و هم سازگار (با قانون و اسپاری گذشته فرانسه).

این سنت بر اساس سه معیار قرار داشته است:

- انتشارات: اسناد گردآوری شده برای یک مخاطب و در خدمت یک هدف عمومی هستند، آن ها نباید در حد حوزه های خصوصی یا ارتباطات درون شرکتی تنزل کنند؛

- رسانه: همه ترتیبات قانون و اسپاری قبلی بر اساس موجودیت فیزیکی یک رسانه قرار داشتند: منابع، چاپی متون موسیقی (1) تصاویر نوارهای صدا دیسکت و جز آن؛ و

- قلمرو: اسناد باید در محدوده مرزهای قلمرو ملی منتشر یا توزیع می شدند.

به طور خلاصه، قانون و اسپاری گذشته برای تمام انتشارات موجود در رسانه های تولید شده یا توزیع شده در قلمرو کشور فرانسه قابل اجرا بوده است.

متأسفانه هیچ یک از این معیارها برای گردآوری [منابع] وبی به راحتی قابل اجرا نیست، زیرا وب گاه ها:

- به ترکیب و ادغام ارتباطات شخصی و عمومی و ترکیب مبتکرانه آن ها تمایل دارند؛

- به معنای دقیقه کلمه، رسانه نیستند، بلکه بیش تر سکویی هستند که سایر رسانه ها تمایل به انتقال به آن ها دارند (می توان کتاب، تصاویر، فیلم، متون موسیقی، و جز آن را در وب یافت)؛ و

- به سهولت، به یک سرزمین خاص - دست کم در مقیاسی وسیع - محدود نمی شوند؛

به علاوه، دیگر نمی توان زبان فرانسه را معیار انتخاب قرار داد زیرا قانون و اسپاری بدون توجه به زبان انتشاراتی اعمال می شود: مجموعه های و اسپاری شده به کتابخانه ملی فرانسه شامل اقلامی به زبان های خارجی می شود که در کشور فرانسه منتشر چاپ یا توزیع شده باشند به همین ترتیب، نمی توانیم از راهبرد قانون و اسپاری خود انتظار تمرکز بر موضوع نویسنده یا سطوح انتشاراتی خاص داشته باشیم جنبه مهم این قانون آن است که مجموعه ها باید منعکس کننده جامعه و فرهنگ فرانسه به هر شکل و صرف نظر از ارزش علمی یا محبوبیت آن ها باشند.

اگر مجبور به انتخاب هستیم منظور بیش تر تهیه نمونه هایی از منابع است تا انتخاب نسل های بعدی تصمیم می گیرند که چه چیزی ارزشمند است نه کتابخانه ملی در قفسه های کتابخانه، نویسندگان ناشناخته و مجله های مستهجن در کنار متفکران بزرگ قرار گرفته اند، انتظار داریم این فلسفه در مورد منابع و بی هم اعمال شود برداشت، فل های فرصت بزرگی را برای عرضه این رویکرد در مقیاس وب فراهم می کند قانون و اسپاری درباره محتوا و قالب - یا محمل - است به این مفهوم که کتابخانه ملی فرانسه به ساخت مجموعه هایی که

منعکس کننده گرایش هایی از نمونه های انتشاراتی است توجه می کند تلاش می کنیم طیف گسترده اشیا و قالب هایی را برای ارائه فراهم کنیم که به طور عملی توسط افراد

ص: 256

scores -1

در نتیجه، برای دامنه ملی خود نیازمند تعریفی بودیم که با وجود انعطاف پذیری و سادگی اجرا منعکس کننده «روح» این گذشته باشد. تعریف یک «کانون» فرانسه در عین داشتن قابلیت انعطاف پذیری، تنها راه تسهیل شیوه های گردآوری اکتشافی در مقیاس کلان بود ثابت شده است، زبان، مکان جغرافیایی اسامی یا موضوع برای تعیین تفاوت در مقیاس کلان بسیار بی ربط و چالش برانگیز هستند، بنابراین، آن چه به عنوان نقطه آغاز احتمالی برای اکتشاف دامنه ملی باقی ماند استفاده از دامنه سطح بالای ملی (تی. ال دی) Top Level domain (1) (TLD) بود این دامنه به سبب تهیه فهرست آغازین اصلی برای اکتشاف با سایر راهبردها تلفیق شده است. از این رو خط مشی واسپاری برخط ما حین انتظار آشکار شدن در حکم صادر شده، به شکل زیر تعریف شد. ما «فرانسه» را چنین در نظر گرفتیم:

- به عنوان یک اصل هر وبگاهی که در دامنه سطح بالای fr. یا هر دامنه سطح بالای مشابه دیگر، با ارجاع به قلمرو رسمی فرانسه ثبت می شود (برای، نمونه re برای جزیره فرانسوی لاری نیو (2))؛

- هر وب گاهی (شاید خارج از فرانسه) که تولید کننده اش در سرزمین فرانسه قرار دارد (معمولاً این مسئله را می توان در وب گاهی یا با استفاده از (3) Domain Name System تعیین کرد)؛

- هر وبگاهی (شاید خارج از فرانسه) که بتواند برای نمایش محتواهای تولید شده در سرزمین فرانسه تأیید شده باشد (این معیار آخر چالش برانگیزتر از آن است که بررسی شود اما مجالی را برای تفسیر و گفت و گو میان کتابخانه ملی فرانسه و تولیدکنندگان وب به وجود می آورد).

البته انتظار نمی رود هیچ یک از این معیارها قبل از این که ما فهرست هسته (4) را بسازیم، به طور کامل تأمین شوند (این امر مانع اکتشاف می شود و در بررسی وب گاه ها، قبل از خزش (5)، ایجاد اخلاص می کند، که به سادگی قابل اندازه گیری نیست). با وجود، این آن ها قصد شان خدمت در چارچوبی قانونی و فرهنگی است به منظور:

- تعریف و توصیف شیوه کلی سیاست آرشیو سازی وب ملی برای مردم صاحبان محتوا، و کتابداران درگیر در پروژه؛

- تبیین وظایف و دستورالعمل های مورد نیاز برای خود تا وقت نظارت بر خزش ها به یاد داشته باشیم؛ تعیین عناصر عینی تصمیم گیری وقتی ما نیازمندیم بدانیم آیا وب گاهی به طور کامل در حوزه کار ماست یا خیر این امر به ویژه زمانی که وب مستر یک وب گاه خاص از کتابخانه ملی می خواهد کار خزش را متوقف کند یا حتی هنگام رویارویی با یک دعوی، قانونی سودمند است: اگر وب گاه در هیچ یک از معیارهای فوق، ننگنجد باید از خزش های آینده و از مجموعه کنار گذاشته شود.

بنابراین، رویکرد دامنه ملی مصالح های را میان قانون واسپاری گذشته و ویژگی های چالش برانگیز

seed list -4 (فهرست یو. آر.ال).

crawl -5

وب نشان می دهد، همچنین مصالح های میان رویکرد کاملاً باز و نامعقولانه ای که به احتمال می توانست ما را طوری هدایت کند تا کل شبکه جهانی وب را به عنوان فرانسه بالقوه در نظر بگیریم، یا بر عکس، رویکردی محدود کننده، که فقط به دامنه سطح بالای fr تنزل می یافت در حالی که معروف است این دامنه تنها بخش محدودی از وب گاه های فرانسوی را شامل می شود در یک کلام، هدف این رویکرد، تأمین تمرکز و انعطاف پذیری است.

2.1. در حال حاضر کجا هستیم

کتابخانه اجازه دارد راه های متفاوتی را برای گردآوری اینترنت فرانسه به کار ببرد: «مؤسسه های قیم ممکن است منابع اینترنت را با به کارگیری فنون خودکار یا با تنظیم موافقت نامه های خاص و فرآیندهای واسپاری با همکاری تولیدکنندگان گردآوری کنند» (ماده 41 II) به همین منظور کتابخانه ملی فرانسه مدلی تلقیقی متشکل از سه راهبردی زیر را تعریف کرده است.

- برداشت فله ای اینترنت فرانسه هدف گردآوری دامنه fr، دست کم به صورت سالانه است. این خزش های فراگیر (1) امکان آرشیو تصاویر وب فرانسه را به کتابخانه می دهند. این رویکرد، در مقایسه با هزینه های ماشینی و نیروی انسانی برداشت [و حجم] اطلاعات بازاریابی، شده باصرفه تر است. با وجود، این به خاطر محدودیت منابع و ابزاری با چنین خزش گر (2) هایی امکان ندارد بتوان وب عمیق (وبگاه های بسیار بزرگ پایگاه های اطلاعاتی و جز آن) را گردآوری کرد.

- خزش های کانونی (3) برای وب گاه ها به تعداد محدود این سایت ها در داخل یا خارج از کتابخانه ملی، فرانسه با همکاری شبکه ای از کتابداران و پژوهش گران کشف شده اند خزش گر های کانونی به وبگاه های بزرگ و به وب گاه های غالباً در حال اصلاح، اختصاص دارند.

- واسپاری های الکترونیکی ویژه شمار محدودی انتشارات الکترونیکی

کتابخانه، منتظر تصویب قانونی که فنون خزش را بررسی کند نشد پروژه آرشیو سازی وب در 1999 آغاز شد. نخستین خزش گر کانونی رویداد محور حدود سال 2002 آزمایش شد در آن زمان کتابخانه نزدیک به دو میلیون وب گاه وابسته به انتخابات کشور (انتخابات ریاست جمهوری و مجلس) را گردآوری کرد این کار برای اروپا و برای انتخابات محلی دو سال تمدید شد کتابخانه ملی فرانسه 1162 وب گاه را گردآوری کرد (4)

با وجود این، تجهیزات فنی (سخت افزاری و نرم افزاری) مهارت ها و تجربه لازم برای تحقق خزش گر های بزرگ مقیاس اینترنت، فرانسه در کتابخانه ملی هنوز کافی نبود این دلیل چرایی موافقت نامه همکاری کتابخانه ملی فرانسه با آرشیو اینترنت (IA)، بنیاد غیرانتفاعی در آرشیو سازی شبکه جهانی وب، از 1996 است. در نوامبر 2004 این دو مؤسسه موافقت نامه تحقیقاتی را با نام «پروژه پژوهشی: انتخاب

ص: 258

broad crawls -1

crawler -2

Focused crawls -3

4- خزش گرها برای انتخابات 2002 و 2004 استفاده شدند برای اطلاعات بیش تر درباره این دو خزش گر به منبع 19 مراجعه کنید.

دامنه ملی برای آرشیوسازی وب)) امضا کردند هدف این پروژه تعیین شیوه ها و ابزارهایی برای استفاده در یک خزش دامنه ملی وب بود موافقت نامه تصریح کرد برای خزش های فراگیر، لازم است بررسی این ابزارها و شیوه ها توسط آرشیو اینترنت (IA) انجام و داده هایی که طی این خزش ها گردآوری می شد به قفسه های ذخیره سازی بی.ان.اف تحویل شود.

نخستین خزش فراگیر در چارچوب این موافقت، نامه در پایان سال 2004 اتفاق افتاد. از آن، در سال های 2005 2006 و 2007 سه خزش فراگیر fir انجام شد (این آخرین خزش که تا سال 2009 ادامه خواهد، داشت مرهون بسط موافقت نامه تحقیقاتی است).

خزش های با مقیاس کوچک تر نیز - مستقیم یا غیر مستقیم - توسط کتابخانه ملی فرانسه اجرا شده. است دو خزش کانونی در شمار محدودی وب گاه (در حدود 4000 عدد) توسط IA برای پروژه تحقیقاتی انجام شد [15] از 2007، این وب گاه ها توسط بی.ان.اف با استفاده از امکانات خود، برداشت می شوند. در همان سال سایر خزش های کانونی موضوعی یا رویداد محور به اجرا درآمدند، مانند وب گاه های مربوط به انتخابات ملی 2007 فرانسه.

به این ترتیب بی.ان.اف تاکنون - مستقیم یا به لطف همکاری با IA- چهار خزش فراگیر به علاوه حجم فراوانی خزش های کانونی به اجرا در آورده است، دیگر آرشیوسازی، وب یک پروژه در کتابخانه ملی فرانسه نیست بلکه فعالیتی روزمره و واحدی دائمی در اداره قانون واسپاری (1) این کتابخانه است. از آن جا که این امر بهترین رویکرد در مواجهه با چالش گردآوری حجم روزافزون اشیای دیجیتال در وب است برداشت فله ای هنوز اولویت نخست به شمار می آید.

با وجود این برداشت فله ای به معنای برداشت کورکورانه نیست حتی وقتی روبات ها سعی دارند حداکثر فایل های وبی را کشف کنند خود را به رعایت قوانین و تنظیمات ملزم می کنند [20]. تصمیمات فنی، قبل حین و بعد از خزش نقش قطعی در نتیجه برداشت دارد.

این مقاله راهبردهایی که بی.ان.اف با همکاری آرشیو اینترنت (IA) برای اجرای خزش های بزرگ مقیاس ارائه کرده است توصیف میکند این خزش ها می تواند با دیدگاه دامنه ملی فرانسه سازگار باشد.

2. طراحی خزش

2.1. هدف چیست؟

به نظر می رسد، اجرای خزش فراگیر در حدود صدها میلیون فایل با طیفی از مشکلات فنی اجتناب ناپذیر همراه باشد. با این حال، نخستین سؤالی که قبل از آغاز یک خزش باید پاسخ داده شود، سؤالی فنی نیست: هدف این خزش گر چیست؟ پاسخ ها بسته به اینکه از یک شرکت بزرگ نهادی پژوهشی، یا مؤسسه ای میراثی باشد متفاوت خواهد بود. اگر لازم است داده ها در مدت طولانی نمایه سازی (توسط موتور جست و جو)، تحلیل (به طور مثال برای شناسایی دامنه های وبی) و یا آرشیو و نمایش داده شوند، خزش فراگیر شیوه یکسانی را اجرا نمی کند.

ص: 259

به عبارت دیگر توجه به محدودیت های فنی، هنگام تعریف اهداف خزش فراگیر ضروری است. این نکته دلیل لزوم گفت و گوی دائمی میان کتابداران (که وظیفه تعریف خط مشی مجموعه سازی را بر عهده دارند) و مهندسان (که وظیفه اجرای خزش ها را بر عهده دارند) را به خوبی بیان می کند متدولوژی ای که در این بخش توصیف می شود نقطه تلاقی دغدغه های مربوط به این دو گروه است

خزش های کتابخانه در چارچوب قانون و اسپاری به اجرا درآمده اند سودمندی این چارچوب تنها برای پرداختن به مسائل حفاظت مالکیت معنوی نیست بلکه آرشیوسازی وب را به عنوان وظیفه ای مستمر مطرح می کند. مطابقت خط مشی مجموعه سازی در آرشیوهای وب، با قالب های قدیمی، انتشارات امری ضروری است.

کشف خودکار وب گاه ها با رویت راه کاری است برای سازگاری با ویژگی «غیر تبعیض آمیزانه» قانون واسپاری، فرانسه تا هم «بهترین» (انتشارات ادبی علمی) و هم «بدترین» (از آگهی ها گرفته تا پورنوگرافی) انتشارات فرانسوی را گردآوری کند با وجود، این حتی رویت ها در معرض پیش داوری هستند. ساختار فرایبندی، وب نخست به کشف و گردآوری پر استناد ترین وب گاه منتهی شد اما با رویت های آرشیوسازی ما وب گاه های کم طرفدار فراموش نمی شوند برای اجتناب از این پیش داوری BNF در سپتامبر 2007، موافقت نامه ای را با انجمن همکاری نامگذاری اینترنت فرانسه (آفنیک) (1)، سازمان مسئول دامنه های fr و re. امضا کرد. طبق شرایط، موافقت نامه این سازمان باید هر دو سال یک بار فهرست کاملی از اسامی دامنه ثبت شده موجود در دامنه های fr و re را ارائه دهد (هم اکنون بیش از یک میلیون اسامی دامنه دارد) به عبارت دیگر BNF باید اعتبار این اطلاعات ارزشمند را تضمین کند.

بنابراین، هدف یک خزش دامنه، بزرگ گردآوری نمونه ای جامع از دامنه ملی و ارائه تصویری از تولیدات [فکری] فرانسه در زمان این گردآوری است در بیش تر موارد نمونه به عنوان یک تصویر در نظر گرفته می شود- راه کاری برای حفظ و تثبیت یک فضای در حال تحول از آن جا که گردآوری همه چیز ممکن نیست به بهای [از دست دادن] یکی برداشت چند سند از هر وب گاه را به گردآوری کل چند وب گاه ترجیح می دهیم

به این سبب خزش های عمیق تری را در مهم ترین وب گاه ها تجربه نکردیم مانند کتابخانه ملی استرالیا که با گذاشتن اولویت بالا در مورد وب گاه های دولتی و دانشگاهی این امر را انجام داد فهرستی از این وب گاه ها توسط کتابداران منتشر شده است [17]) این کار لازم هم نبود زیرا خزش های فراگیر ما با همراهی خزش های کانونی به صورت کامل به اجرا در می آیند.

از سوی دیگر، احتمال تأثیر سریع وب از تحولات فناورانه وجود دارد. قالب های انتشاراتی نوین پدیدار می شوند و در عرض چند ماه گسترش می یابند خزش فراگیر باید بازتاب دهنده این تحولات باشد به عنوان نهاد قانون واسپاری یکی از اهداف، ما روشن کردن قالب های نوین انتشاراتی است و در نتیجه کسب اطمینان از اینکه رویت توانایی برداشت این اسناد را داشته باشد به این سبب، در سال 2006 ، به وب نوشت ها و وب گاه های شخصی توجه بیش تری داشتیم و در سال 2007 بر ویدیو ها تأکید

ص: 260

داریم (در ادامه می آید).

طبق اهداف از پیش تعریف شده در آینده احتمال دارد مجموعه قبل و در زمان خزش شکل گیرد. قبل از آغاز کار، برداشت درباره دو عامل مهم که سهم بزرگی در ساخت مجموعه های آینده دارند تصمیم گیری شد طراحی فهرست هسته، و تنظیمات خزش.

2.2 فهرست هسته

خزش گر وظایفش را با فهرستی از یو. آر. ال آغاز می کند که فهرست هسته نامیده می شود. هسته ها، در هایی برای دستیابی به وب گاه ها به شمار می آیند؛ به این سبب کیفیت برداشت تا حد زیادی به کیفیت این فهرست بستگی دارد.

از آن جا که همکار ما در اجرای خزش فراگیر برای چهار سال همان [آرشیو اینترنت (IA)] باقی ماند، امکان غنی کردن تدریجی فهرست هسته وجود داشت منابع مختلف هسته ها سال به سال افزایش می یافتند:

- 2004: فهرست هسته [یو.آر.ال]، برای نخستین خزش فراگیر از استخراج دامنه های fr. آخرین خزش الکسا به وجود آمد. (1)

- 2005: هسته های حاصل از استخراج دامنه های fr آخرین خزش گر و میزبان آلکسا، هنگام خزش فراگیر قبلی کتابخانه ملی فرانسه کشف شدند هدف فرصت دادن به خزش گر برای پیش رفتن و کشف میزبان های جدید بود

- 2006: فهرست هسته به شیوه سال گذشته آغاز شد.

- 2007: به لطف امضای توافق نامه با آفنیک فهرست جامع اسامی دامنه fr و re. به عنوان فهرست هسته به کار رفت برای اطمینان از تداوم و سازگاری با خزش های فراگیر، قبلی این فهرست با استخراج هایی از خزش آلکسا و از خزش های میزبان قبلی ادغام شد.

آفنیک شامل:

- 890064 دامنه fr

- 1516 دامنه re، و

- 21 / 753 اسم دامنه سطح دوم دامنه سطح دوم (2،3،3) را ببیند.

با وجود این استفاده از فهرست، آفنیک به عنوان فهرست هسته، ساده امکان پذیر نبود این فهرست شامل اسامی دامنه می شد نه یو آر ال آن ها به عبارتی پشت نام هر دامنه یک وب گاه وجود نداشت

بنابراین در برخورد با فهرست، آفنیک چند تحلیل مورد پردازش قرار گرفتند.

ص: 261

هایی درباره سایت های مورد بررسی فراهم کند و بر اساس داده های گردآمده از دیگر کاربران ، صفحه های مرتبطی را که ممکن است مورد علاقه آن ها باشند، توصیه می کند این شرکت از ، 1996 آرشیوهای خزش خود را در اختیار آرشیوهای اینترنتی قرار داده است [16]

نخستین آن‌ها کمی کردن تعداد دامنه‌هایی بود که هنوز فعال بودند و باید به طور عملی، به عنوان هسته، مورد استفاده قرار می‌گرفتند. این بررسی سنگین توسط آرشیو اینترنت (IA) انجام شده است. برای کسب اطمینان از برخط بودن پاسخ آن‌ها هر دامنه با دو نشانی مختلف بررسی می‌شد (یعنی بررسی 1/780/128 یو.آر.ال):

<http://domainname.fr> و <http://www.domainname.fr>.

عکس

۲۶۲ مدیریت منابع اطلاعاتی وب

نخستین آن‌ها، کمی کردن تعداد دامنه‌هایی بود که هنوز فعال بودند و باید به طور عملی، به عنوان هسته، مورد استفاده قرار می‌گرفتند. این بررسی سنگین، توسط آرشیو اینترنت (IA) انجام شده است. برای کسب اطمینان از برخط بودن پاسخ آن‌ها، هر دامنه با دو نشانی مختلف بررسی می‌شد (یعنی بررسی ۱/۷۸۰/۱۲۸ یو.آر.ال):

<http://domainname.fr> و <http://www.domainname.fr>.

هر دو صورت دامنه دارای دی.ان.اس است	۷۹ درصد
یک صورت دامنه دارای دی.ان.اس است	۱۴ درصد
هیچ یک از دو صورت دی.ان.اس ندارد	۷ درصد

شکل ۱. پاسخ به وجود دی.ان.اس^۱ در اسامی دامنه آفتیک

بر اساس نتایج، فهرستی از یو.آر.ال‌های معتبر به دست آمد، تا آنجا که، اگر نسخه WWW پاسخ نمی‌داد، از «domainname.fr»، یا <http://www.domainname.fr> یا <http://domainname.fr> استفاده می‌شد. هسته‌های بدون دی.ان.اس در فهرست هسته به طور تصادفی وارد می‌شدند.

ما اسامی دامنه‌های (و به‌ویژه در مجموعه خزش فراگیر ۲۰۰۶) موجود در فهرست آفتیک را که در آرشیوهای وب نیز موجود بودند، بررسی کردیم. فهمیدن این نکته که فقط ۳۰ درصد از دامنه‌های آفتیک در مجموعه ما وجود داشت، بسیار حائز توجه بود.

این امر شاید به افزایش چشمگیر اندازه دامنه .fr مربوط باشد. به یمن شدن پی در پی قوانین مختص .fr، از سال ۲۰۰۴، به‌طور مداوم افزایش یافته است. عمده‌ترین آنها، اجازه ایجاد دامنه سطح بالا برای اشخاص، در سال ۲۰۰۶ بود (تا این تاریخ، تنها اداره‌ها و انجمن‌های خصوصی مجوز ایجاد دامنه .fr را داشتند): با گذشت یکسال، .fr، بیش از ۶۳ درصد افزایش یافت. این امر باید مربوط به ایجاد تصور مثبت از ccTLD در کاربران اینترنت باشد. اشخاص ۳۰ درصد دامنه‌های .fr ثبت شده، و ۵۰ درصد ثبت‌نام‌های جدید را تشکیل می‌دهند [۲].

همچنین، وجود پدیده فوق، در فهرست آفتیک، ممکن است تاحدی به خاطر حضور اسامی دامنه-هایی باشد که یا پیوندی به وبگاه‌های دیگر ندارند و یا پیوندشان ضعیف است، و توسط رویات‌ها در سال‌های گذشته کشف نشده بودند.

این رشد فوق‌العاده، تا حد زیادی گسترش چشمگیر فهرست هسته را از ۲۰۰۶ تا ۲۰۰۷ توصیف می‌کند.

۱. DNS [سروری که دامنه‌های آدرس سایت‌ها را پشتیبانی می‌کند، مترجم].

شکل 1. پاسخ به وجود دی. ان. اس (1) در اسامی دامنه آفنیک

بر اساس نتایج، فهرستی از یو. آر. ال های معتبر به دست آمد تا آن جا که اگر نسخه www پاسخ نمی داد، از «domainname.fr» یا http://www.domainname.fr یا http://domainname.fr استفاده می شد.

هسته های بدون دی. ان. اس در فهرست هسته به طور تصادفی وارد می شدند.

ما اسامی دامنه های (و به ویژه در مجموعه خزش فراگیر 2006) موجود در فهرست آفنیک را که در آرشیوهای وب نیز موجود بودند بررسی کردیم فهمیدن این نکته که فقط 30 درصد از دامنه های آفنیک در مجموعه ما وجود داشت بسیار حائز توجه بود

این امر شاید به افزایش چشمگیر اندازه دامنه .fr مربوط باشد به یمن ساده شدن پی در پی قوانین مختص .fr. از سال 2004 .fr به طور مداوم افزایش یافته است. عمده ترین، آن ها اجازه ایجاد دامنه سطح، بالا برای اشخاص در سال 2006 بود (تا این، تاریخ تنها اداره ها و انجمن های خصوصی مجوز ایجاد دامنه .fr را داشتند): با گذشت یکسال .fr. بیش از 63 درصد افزایش یافت این امر باید مربوط به ایجاد تصور مثبت از ccTLD در کاربران اینترنت باشد. اشخاص 30 درصد دامنه های .fr ثبت شده و 50 درصد ثبت نام های جدید را تشکیل می دهند [2].

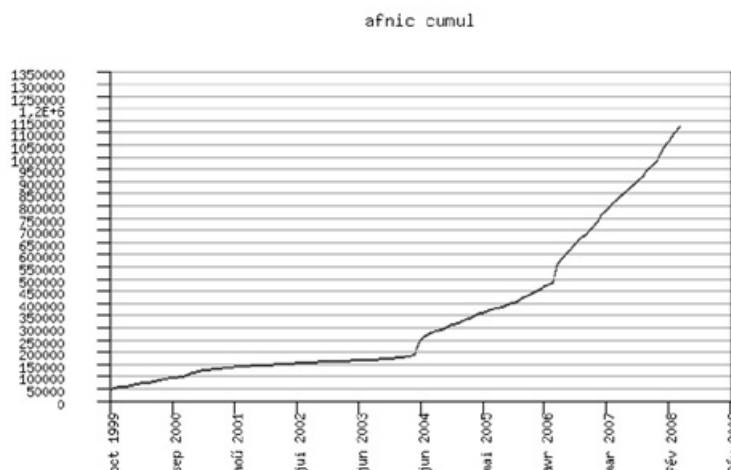
همچنین، وجود پدیده فوق در فهرست، آفنیک ممکن است تا حدی به خاطر حضور اسامی دامنه هایی باشد که یا پیوندی به وب گاه های دیگر ندارند و یا پیوندشان ضعیف است و توسط روایات ها در سال های گذشته کشف نشده بودند.

این رشد فوق العاده تا حد زیادی گسترش چشم گیر فهرست هسته را از 2006 تا 2007 توصیف می کند.

ص: 262

1-DNS [سروری که دامنه های آدرس سایت ها را پشتیبانی می کند مترجم].

قانون واسپاری وب فرانسه... ۲۶۳



شکل ۲. تغییر حجم دامنه .fr از ۲۰۰۶ تا ۲۰۰۷

۲۰۰۷	۲۰۰۶	
۵۸/۲۲۴	۲۰۷/۰۴۶	هسته‌های استخراج شده از خزش‌های آلکسا
۲/۲۹۵/۸۹۰	۴۲۷/۴۷۶	هسته‌های استخراج شده از خزش‌های قبلی
۸۹۰/۰۶۴	-	فهرست آف‌نیک
۲/۸۸۸/۷۲۳	۵۶۲/۶۳۷	تعداد کل پس از کاهش تکراری‌ها

شکل ۳. تغییر اندازه فهرست هسته از ۲۰۰۶ تا ۲۰۰۷

۳.۲. تنظیمات خزشگر

برداشت‌ها توسط هریتریکس^۱ انجام شده بود؛ خزشگر معتبر آرشیوی وب منبع باز توسط آرشیو اینترنت (IA) با کمک‌هایی از اعضای کنسرسیوم آی آی پی سی^۲ (به‌ویژه کتابخانه‌های ملی کشورهای حوزه شمال اروپا)^۳ [۱۱ و ۲۱]. این خزشگر، در کتابخانه ملی فرانسه برای خزش‌های فراگیر، از ۲۰۰۴ و برای خزش‌های کانونی از ۲۰۰۶، مورد استفاده قرار گرفته است.

۱. قابل دسترس در: <http://www.afnic.fr/actu/stats/evolution> (دیده شده: ۱۶ مارس ۲۰۰۸).

۲. Heritrix

۳. IIPC

۴. Nordic national libraries

۵. کنسرسیوم بین‌المللی حفاظت اینترنت در ۲۰۰۳ توسط ۱۲ سازمان (آرشیو اینترنت و چند کتابخانه ملی) برای یافتن راه حل‌های مشترک به منظور آرشیوسازی و اطمینان از دستیابی بلندمدت به انتشارات الکترونیکی در وب تأسیس شد. از ۲۰۰۷، عضویت برای سازمان‌های جدید باز است. برای اطلاعات بیشتر در باره اهداف و فعالیت‌های کنسرسیوم منبع [۱۳] را ببینید.

afnic cumu1

شکل ۲. تغییر حجم دامنه .fr از ۲۰۰۶ تا ۲۰۰۷ (۱)

شکل ۳. تغییر اندازه فهرست هسته از ۲۰۰۶ تا ۲۰۰۷

برداشت ها توسط هریتریکس (2) انجام شده بود؛ خزش گر معتبر آرشیوی وب منبع باز توسط آرشیو اینترنت (IA) با کمک هایی از اعضای کنسرسیوم آی. آی. پی. سی (3) (به ویژه کتابخانه های ملی کشورهای حوزه شمال اروپا) (4) [11 و 21] (5) این خزش گر در کتابخانه ملی فرانسه برای خزش های فراگیر از 2004 و برای خزش های کانونی از 2006 مورد استفاده قرار گرفته است.

ص: 263

1- قابل دسترس در <http://www.afnic.fr/actu/stats/evolution> (دیده شده: 16 مارس 2008)

Heritrix -2

IIPC -3

Nordic national libraries -4

5- کنسرسیوم بین المللی حفاظت اینترنت در 2003 توسط 12 سازمان (آرشیو اینترنت و چند کتابخانه ملی) برای یافتن راه حل های مشترک به منظور آرشیوسای و اطمینان از دستیابی بلند مدت به انتشارات الکترونیکی در وب تأسیس شد از 2007 عضویت برای سازمان های جدید باز است. برای اطلاعات بیش تر درباره اهداف و فعالیت های کنسرسیوم منبع [13] را ببینید.

با قابلیت تنظیم، بالا هریتريکس تغيير مقادير زيادي از تنظيمات شامل، دامنه، اولويت هاي خزش، فيلترها، آداب روياتها و جزآن را مقدور مي سازد.

تنظيمات فوق به خزشگر فرمان مي دهد چه چيزي را و چگونه بايد برداشت کند. زيرا تنظيمات تأثير زيادي در مجموعه سازي دارند و رويات لازم است مطابق با اهداف خزش گر قالب بندي شود.

1.3.2. حوزه

حوزه تعيين شده براي خزشگر تعريف مي کند که کدام یک از يو. آر. ال هاي کشف شده بايد در برداشت وارد و کدام یک کنار گذاشته شود.

از اين رو، حوزه خزش کتابخانه ملي، فرانسه شامل هر وب گاهي مي شود که در هر نام دامنه متعلق باشد به:

- دامنه سطح بالاي fr؛

- دامنه سطح بالاي re؛ و

- هر دامنه ديگري که خزش گر با آن مواجه شود زيرا از يك دامنه fr يا از يك re تغيير مسير داده است (طی بررسی خزش، 398,548 تغيير مسير از اسامي دامنه آفنيک مورد توجه قرار گرفتند). به سبب تغيير مسير از <http://yahoo.fr> به <http://fr.yahoo.com> لازم است خزش گر هر چيزي را که بخشي از <http://fr.yahoo.com> است گزينش کند اگر چه در چنين موردی، خزش گر بايد در همان ميزبان باقي بماند به اين ترتيب خزش گر هر چيزي را که بخشي از <http://fr.yahoo.com> است، بر مي گزيند، نه از <http://de.yahoo.com> يا حتي <http://fr.news.yahoo.com>.

از يك سو اين حوزه بسيار بزرگ به نظر مي رسد - شامل بيش از يك ميليون دامنه - و از ديگر سو محدود است زيرا وب فرانسه تنها در دامنه هاي fir ميزباني نمي شود به نقل از گزارش منتشر شده، آفنيک کمتر از 30 درصد وب گاه هاي فرانسه در fir ميزباني مي شود [2] اين رقم با تحليل مجموعه سايت هاي انتخابات، 2007 مورد تأييد قرار گرفته است. تنها 36 درصد يو. آر. ال هاي اين مجموعه براي ميزباني در دامنه fr. پذيرفته شده اند حتي اگر از قبل چنين تعيين کنيم که وب گاه هاي سياسي نماينده تام الاختيار کل وب فرانسه نيستند به طور مثال کاربرد وسيع .org. کل توسط وب گاه هاي احزاب و اتحاديه هاي تجاري به بازنمون بيش از حد اين دامنه سطح بالا منجر مي شود(بايد بدانيم fir در بخش عمده وب گاه هاي فرانسه استفاده نمي شود).

درصد	URLs	TLD
36.11	22938947	.fr
28.93	18373574	.com
25.12	15955225	.org
5.81	3690655	.net
1.39	882656	.de
1.16	733634	.info
257769	257769	.eu
0.17	110944	.tv
0.17	106753	.us
0.16	99012	.re
0.58	368148	سایر TLDs

شکل ۴. تعداد یو.آر.ال‌های هر دامنه سطح بالا (TLD)، خزش‌های کانونی انتخابات ۲۰۰۷

با وجود این، هدف ما برداشت دامنه‌های دیگری غیر از .fr و .re از طریق دنبال کردن پیوندهای تغییر مسیر داده شده بود. این موضوع، راهکاری برای اختیار کردن حوزه‌ای تغییرپذیر است. شاید رویکرد اکتشافی‌تر به گردآوری وبگاه‌های خارجی‌تر منتهی شود: این امر مشکلات قانونی (کتابخانه ملی فرانسه را به خارج از چارچوب قانون واسپاری می‌برد) و اقتصادی را پدید می‌آورد: گردآوری وبگاه‌های غیرمرتبط (مطابق با مأموریت ما) باعث اشغال فضایی می‌شود که از آن می‌توان برای برداشت سایت‌های معتبر فرانسه استفاده بهتری کرد. بنابراین، تأکید بر .fr یک انتخاب واقع‌بینانه بود. به علاوه، از آنجا که .fr به‌طور چشمگیری در حال رشد است، می‌توان امیدوار بود که این دامنه، بخش بزرگ‌تری از وب فرانسه را سال به سال نمایندگی کند. از دیگر سو، ممکن است در آینده، به‌طور مثال با استفاده از مراجعه به دی.ان.اس خودکار (محل جغرافیایی)، راه‌های نوینی را برای کشف سایت‌های فرانسوی خارج از دامنه .fr بررسی کنیم.

۲.۳.۲. اولویت‌های خزش

هر URL که معلوم شود در حوزه خزش است، در «صف» خزشگر قرار می‌گیرد، یعنی در فهرست فایل‌های در انتظار خزش از آنجا که ممکن است روبات نتواند همه یو.آر.ال‌هایی که زمان اجرای یک خزش فراگیر در سر راهش می‌یابد، برداشت کند، برای خزشگر، مدیریت این صف و تعیین اولویت‌های خزش شدن مسائل اساسی هستند.

نخستین تصمیم مهم، انتخاب میان رویکرد «به ازای هر دامنه»^۱ و «به ازای هر میزبان»^۲ بود. با رویکرد به ازای هر میزبان، که در خزش‌های فراگیر قبلی استفاده شده بود، URLها در صف در حال انتظار خزش شدن در هر میزبان دسته‌بندی می‌شدند، و هر میزبان به‌صورت جداگانه تلقی می‌شد. این دسته‌بندی، به

1. per-domain
2. per-host

شکل 4. تعداد یو.آر.ال‌های هر دامنه سطح بالا (TLD) خزش‌های کانونی انتخابات 2007

با وجود این هدف ما برداشت دامنه‌های دیگری غیر از .fr و .re از طریق دنبال کردن پیوند های تغییر مسیر داده شده بود این موضوع راه کاری برای اختیار کردن حوزه ای تغییر پذیر است. شاید رویکرد اکتشافی تر به گردآوری وب گاه های خارجی تر منتهی شود: این امر مشکلات قانونی (کتابخانه ملی فرانسه را به خارج از چارچوب قانون واسپاری می برد) و اقتصادی را پدید می آورد گردآوری وب گاه های

غیر مرتبط (مطابق با مأموریت ما) باعث اشغال فضایی می شود که از آن می توان برای برداشت سایت های معتبر فرانسه استفاده بهتری کرد بنابراین تأکید بر fr یک انتخاب واقع بینانه بود. به علاوه، از آن جا که .fr به طور چشم گیری در حال رشد است می توان امیدوار بود که این دامنه، بخش بزرگ تری از وب فرانسه را سال به سال نمایندگی کند. از دیگر سو ممکن است در آینده به طور مثال با استفاده از مراجعه به دی.ان.اس خودکار (محل جغرافیایی) راه های نوینی را برای کشف سایت های فرانسوی خارج از دامنه fir بررسی کنیم

2.3.2. اولویت های خزش

هر URL که معلوم شود در حوزه خزش است، در «صف» خزشگر قرار می گیرد، یعنی در فهرست فایل های در انتظار خزش از آن جا که ممکن است رویات نتواند همه یو. آر. ال هایی که زمان اجرای یک خزش فراگیر در سر راهش می یابد برداشت کند برای خزش گر مدیریت این صف و تعیین اولویت های خزش شدن مسائل اساسی هستند.

نخستین تصمیم مهم انتخاب میان رویکرد «به ازای هر دامنه» (1) و «به ازای هر میزبان» (2) بود. با رویکرد به ازای هر میزبان که در خزش های فراگیر قبلی استفاده شده بود URL ها در صف در حال انتظار خزش شدن در هر میزبان دسته بندی می شدند و هر میزبان به صورت جداگانه تلقی می شد این دسته بندی به

ص: 265

per-domain -1

per-host -2

خزش وب گاه های بیش تری با چندین میزبان منجر می شد. وب نوشت های میزبانی شده در سکوها های تجاری یا صفحه های، شخصی به عنوان میزبان های متمایز مقادیر زیادی از فضای مجموعه را اشغال می کنند به این ترتیب با استفاده از رویکرد به ازای هر دامنه با این وب گاه ها به عنوان هویتی مستقل برخورد می شود.

برای خزش سال 2007 تصمیم گرفتیم رویکرد به ازای هر دامنه را اتخاذ کنیم. علت اصلی این امر تبعیت از فهرست آف نیک، بود که فقط بر دامنه ها نظارت می. کند همچنین به عنوان یک نهاد و اسپاری می خواستیم به هر دامنه «شانس / فرصت» یکسانی برای برداشت شدن بدهیم. به علاوه، از آن جا که فهرست هسته 2007 بزرگ تر از فهرست های قبلی بود و مقادیر داده هایی که باید توسط آرشیو اینترنت (IA) بازیابی می شد افزایش نمی یافت (قسمت 2، 5 را ببینید) بیم داشتیم، که عمق خزش کم تری را کسب کنیم. بنابراین به وب گاه های تجاری یا سازمانی که اغلب چندین میزبان دارند و با خزش های کانونی بهتر برداشت می شوند، اهمیت زیادی بدهیم

در آغاز خزش لازم بود اسامی بزرگی مانند free.fr skyblog.fr orange.fr همچون سایر دامنه ها مدیریت شوند. تصمیم گیری برای اقدام به رفتار خاص با این وب گاه ها به تعویق افتاد: در طول خزش، هنگام تحلیل گزارش های ارسالی توسط مهندسان آرشیو اینترنت (IA) امکان انتخاب های مرتبط تر به نظر آسان تر می نمود.

همچنین برای جلوگیری از نمایندگی بیش از حد وب گاه های بزرگ برای هر سایت بیش ترین «بودجه» را تعیین کردیم برای روبات از یک دامنه برداشت بیش از 10,000 URL ممنوع بود. معنای این، محدودیت این نبود که خزش گر وقتی به این مقدار می رسید به طور کامل متوقف شود؛ اگر کل بودجه خزش صرف نشده بود روبات اجازه داشت مقدار بیش تری برداشت کند

سرانجام برای اطمینان از اینکه روبات منابع کافی برای برداشت کل فهرست هسته در اختیار داشته، باشد، سطح «میزان دوباره پر کردن» (1) را پایین تر انتخاب کردیم زمان اتصال به هسته هر «ریسمان» روبات (2) دستور برداشت صد URL که با دامنه مشابه را دریافت می کرد و سپس به هسته دیگر می رفت. پس از اتمام برداشت صد URL نخست هر هسته روبات به همان هسته باز می گشت.

این، رویکرد متوسطی میان رویکرد نخست پهنا (3) با هدف گردآوری وب گاه های مختلف تا جای ممکن و رویکرد آرشیو سازی ناب (4) است که در آن خزش عمیق یک سایت قبل از [آغاز] برداشت سایت بعدی انجام می گیرد

تصمیم گیری های فوق خطری را برای کتابخانه به همراه خواهد داشت در حقیقت، مجموعه فراهم شده نسبت به مجموعه های سال های قبل بسیار متفاوت است. با وجود این، برای ما بسیار اهمیت داشت که برای کشف راه های جدید در برآوردن ضرورت های وظیفه واسپاری قانونی خود، این شیوه گردآوری را تجربه کنیم.

ص: 266

replenishamount -1

thread" of the robot" -2

breadth-first -3

pure archiving approach -4

تقاضا رویکرد در هر دامنه این بود که: برای روایات ممکن نبود به طور خودکار دامنه های سطح دوم (اس.ال.دی) (1) را کشف کند در مورد کار ما دامنه های سطح دوم بخش های فرعی تخصصی fr هستند. این دامنه ها یا «اسامی بخش سطح دوم» (2) مختص شناسایی یک صنعت یا یک بخش متعارف (3) مانند aeroport.fr برای فرودگاه ها یا gov.fr برای وب گاه های دولتی هستند یا «دامنه های توصیف سطح دوم» (4) که برای شناسایی یک فعالیت یا عنوان چند نوع، اختصاص داده شده اند (به طور مثال، asso.fr. برای فدراسیون ها یا tm.fr برای دارندگان نشان های تجاری)

بدون هیچ تنظیم خاصی وب گاه های متفاوت که در دامنه های سطح دوم مشابه میزبانی شده اند، از نظر روایات ها به عنوان یک هویت محسوب شده و بودجه یکسانی را دریافت نمی کنند برای اجتناب از این مسئله، فهرست کلی دامنه های میزبانی شده در دامنه های سطح دوم را با دستور برخورد با این وب گاه ها به عنوان وب گاه های منفرد به روایات دادیم (این فهرست توسط آفیک نیز برای کتابخانه ملی فرانسه تهیه شده است).

به محتوایی که در صفحه ای جای داده شده بود اما در دامنه ای متفاوت از خود آن صفحه میزبانی می شد توجه خاصی مبذول شد به روایات اجازه داده شد از سه جهش انتقالی حداکثر (5) پیروی کند که پیوندهایی جاسازی شده در صفحه های وبی هستند. انجام این امر برای برداشت تعدادی زیادی فایل، ویدئویی که هدف مهم این خزش بود ضرورت داشت.

در واقع ما سعی داشتیم راه حل های بیشتری را برای برداشت فایل های ویدئویی بیابیم. تعداد ویدئو در وب به طور مداوم رو به افزایش است و این مسئله برای روایات هایی که آن ها را گردآوری می کنند چالش برانگیز است. مشکلات برداشت فایل های ویدئویی به اندازه و شیوه پخش آن ها مربوط می شود (برای تحلیل مشکلات گردآوری رسانه های کنونی، منابع [4] و [7] را ببینید). کمبود فایل های ویدئویی، به سبب وجود محدودیت های فنی در خزش گر، با هدف ما از آرشیو سازی یعنی تهیه تصویری «بازنمون» از وب تناقض دارد ما تصمیم گرفتیم تلاش های خود را بر دو سکوی اصلی پخش ویدئو مورد استفاده netsurfers (6) بر اینترنت فرانسه متمرکز سازیم یوتوب (7) و معادل فرانسوی آن دیلموشن (8) Dailymotion ما از مهندسان آرشیو اینترنت خواستیم مهندسی فنی دیلموشن را مطالعه کنند (آن ها قبلا یوتوب را به شکل بسیار خوبی بررسی کرده بودند) و نسخه های هریتریکس بهبود بخشند تا خزش گرها امکان گردآوری ویدئو فایل های این وب گاه ها را داشته باشند.

ص: 267

1- Second Level Domains (SLD)

2- second level sector names

3- regulated sector

4- second level descriptive domains

5- max-transshops

6- موجسواران

7- [25]

۲۰۰۷	۲۰۰۶	۲۰۰۵	۲۰۰۴	
دامنه	میزبان	میزبان	میزبان	قلمرو
۱۰۰	۱۰۰	۱۰۰	---	جهش - حداکثر
۳	۳	۳	---	جهش انتقالی حداکثر
۱۰۰	۵۰۰	۱۰۰۰	---	مقدار بازذخیره (در URLها)
۱۰۰۰۰	۲۰۰۰۰۰	۲۰۰۰۰۰	---	بودجه (در URL)
۴	۴	۵	---	عامل تأخیر
۵۰۰۰	۵۰۰۰	۵۰۰۰	---	حداقل تأخیر (میان دو درخواست برای یک میزبان، به میلی ثانیه)
۱۰۰۰۰	۱۰۰۰۰	۱۰۰۰۰	---	حداکثر تأخیر (میان دو درخواست برای یک میزبان، به میلی ثانیه)
مناسب	مناسب	مناسب	مناسب	تغییر مسیر
100 Mo	100 Mo	100 Mo	---	حداکثر اندازه URLهای بارگذاری شده

شکل ۵. تنظیمات خزشگر، از ۲۰۰۴ تا ۲۰۰۷

۲.۴. پروتکل حذف روبات‌ها

قانون حق مؤلف ۲۰۰۶، به کتابخانه ملی فرانسه، اجازه سرپیچی از پروتکل حذف روبات‌ها (آر.ای.پی)^۱ را می‌دهد: «تولیدکنندگان یا ناشران] نباید در وبگاه‌های خود، رمزگذاری کنند یا مانع دستیابی مؤسسه‌های مسئول برای برداشت شوند^۲». به همین سبب، این کتابخانه، هنگام اجرای خزش‌های کانونی داخلی^۳، معمولاً به robots.txt توجه نمی‌کرد. در واقع، اغلب، حذف روبات‌ها برای جلوگیری از خزشگرهای گردآورنده صفحه‌هایی که قرار نیست نمایه‌سازی شود، توسط مدیران وب به‌کار می‌رود: راهنماهای تصویری یا صفحه‌های CSS. اما، این اسناد می‌توانند برای خزشگرهای آرشیوسازی مهم باشند، زیرا برای نمایش شکل اصلی صفحه‌های آرشیوی وب در آینده ضرورت دارند.

با وجود این، کتابخانه تصمیم گرفت از قوانین robots.txt برای خزش fr. سال ۲۰۰۷ خود اطاعت کند، همانطور که برای خزش‌های فراگیر قبلی نیز انجام می‌داد. تجارب خزش کانونی نشان داده بود که مدیران وب وقتی می‌فهمیدند که یک روبات، با شکستن قوانین robots.txt مربوط، مشغول خزش در سایت آنهاست، خوششان نمی‌آمد. به‌طور مثال، یک بلاگر فرانسوی که مورد خزش کتابخانه ملی فرانسه قرار گرفته بود، نامه‌های الکترونیکی خشونت‌آمیزی به کتابخانه فرستاد و یک تله - خزشگر برای پایین آوردن سرعت روبات کتابخانه ایجاد کرد، و پیامی در پست وبلاگش گذاشت تا دیگر مدیران وب را

1. robots exclusion protocol (REP)

۲. همچنین این جمله به این معناست که بی‌ان‌اف، به‌طور مثال زمان دستیابی به داده‌هایی که در وبگاه رایگان نیست، به صورت قانونی حق درخواست رمز عبور و کدهای لازم را از مالک وبگاه برای خزش سایت او دارد.

۳ in-house . یا درون سازمانی

شکل ۵. تنظیمات خزش گر، از ۲۰۰۴ تا ۲۰۰۷

۲.۴. پروتکل حذف روبات‌ها

قانون حق مؤلف ۲۰۰۶، به کتابخانه ملی فرانسه اجازه سرپیچی از پروتکل حذف روبات‌ها (آر.ای.پی) (1) را می‌دهد «تولیدکنندگان یا

ناشران] نباید در وب گاه های خود رمزگذاری کنند یا مانع دستیابی مؤسسه های مسئول برای برداشت شوند». (2) به همین سبب این کتابخانه هنگام اجرای خزش های کانونی داخلی (3)، معمولاً به robots.txt توجه نمی کرد. در واقع، اغلب، حذف روبات ها برای جلوگیری از خزش گر های گردآورنده صفحه هایی که قرار نیست نمایه سازی شود، توسط مدیران وب به کار می رود: راهنماهای تصویری یا صفحه های CSS، اما این اسناد می توانند برای خزش گر های آرشیو سازی مهم باشند زیرا برای نمایش شکل اصلی صفحه های آرشیوی وب در آینده ضرورت دارند.

با وجود این کتابخانه تصمیم گرفت از قوانین robots.txt برای خزش fr. سال 2007 خود اطاعت کند، همان طور که برای خزش های فراگیر قبلی نیز انجام می داد. تجارب خزش کانونی نشان داده بود که مدیران وب وقتی می فهمیدند که یک روبات با شکستن قوانین robots.txt مربوط، مشغول خزش در سایت آن هاست خوش شان نمی آمد به طور مثال یک بلاگر فرانسوی که مورد خزش کتابخانه ملی فرانسه قرار گرفته بود نامه های الکترونیکی خشونت آمیزی به کتابخانه فرستاد و یک تله - خزش گر برای پایین آوردن سرعت روبات کتابخانه ایجاد کرد و پیامی در پست وبلاگش گذاشت تا دیگر مدیران وب را

ص: 268

1- robots exclusion protocol (REP)

2- همچنین این جمله به این معناست که بی ان اف، به طور مثال زمان دستیابی به داده هایی که در وبگاه رایگان نیست به صورت کانونی حق درخواست رمز عبور و کدهای لازم را از مالک وبگاه برای خزش سایت او دارد.

3- in-house . یا درون سازمانی

به گذاشتن تله - خزش گر تشویق کند (1) در واقع صفحه های و بی محدود شده توسط فایل های robots.txt اغلب تله های خزش گر را در خود جای می دهند به علاوه، گاهی هدف قوانین حذف جلوگیری خزش گرها برای جست و جوی URL های است که می توانستند موجب اضافه بار در وب گاه ها شوند (به طور، مثال پیشنهادهایی به یک گروه بحث (2)). وقتی امکان نظارت بر برداشت هر وب گاه شخصی وجود داشته، باشد مسائل فنی و مشکلات تولیدکنندگان وبگاه به راحتی در طول خزش کانونی مدیریت خواهند شد. اما کتابخانه تمایلی برای مدیریت آن ها در مقیاس کلان نداشت و نمی خواست مهندسان آرشیو اینترنت زیر فشار مخالفت های مدیران وب عصبانی غرق شوند (به هر حال، سیاست IA احترام به حذف های robots.txt است).

اغلب مؤسسه های مجری خزش های فراگیر در وب به ویژه برداشت دامنه ملی سیاست احترام به حذف روایات ها را انتخاب کردند. [17]، Nearchive.dk [10] مرکز مجازی (کتابخانه سلطنتی و کتابخانه ایالتی و دانشگاهی) با وظیفه آرشیوسازی دامنه دانمارک ترجیح می دهد از robots.txt چشم پوشی کند، زیرا این روایات بیش تر برای سایت های واقعا مهم شبکه مورد استفاده قرار می گیرد [3] (3) در واقع وب گاه ها با میزبانی محتوای بسیار ارزشمند مانند روزنامه ها یا - سایت های احزاب سیاسی رایج ترین وب گاه هایی هستند که از محدودیت های robots.txt استفاده می کنند [24].

5.2. برنامه ریزی با همکاری آرشیو اینترنت (IA)

برای خزش گر فراگیر سال 2007 آرشیو اینترنت و کتابخانه ملی، فرانسه بر اندازه تعیین شده مجموعه توافق دارند به نظر می رسید حدود URL 300,000,000 برای هماهنگ کردن نیازهای خزش منطقی باشد. در صورت نیاز IA می توانست برای 10-15 درصد افزایش تصمیم بگیرد.

هر دو نهاد با سازماندهی روز به روز خزش موافق هستند حوزه و تنظیمات اصلی برای برداشت باید میان IA و کتابخانه ملی فرانسه به بحث گذاشته شود و کتابخانه درباره آن ها تصمیم گیری کند. مهندسان IA لازم است بر کار ماشین های خزش نظارت داشته باشند و هفته ای دو بار گزارش هایی به کتابخانه ملی ارسال کنند (در ادامه می آید) برداشت باید قبل از پایان سال خاتمه یابد و تا ماه های نخست سال 2008، داده ها (برای دسترسی ماشین های Wayback و NutchWAX) نمایه سازی شوند؛ و پس از آن، آرشیوها به قفسه های ذخیره خود در کتابخانه ارسال شوند (4) دو مهندس IA باید برای کمک به کتابخانه برای نصب، قفسه ها مشاوره به تیم BNF و اطمینان از کیفیت مجموعه به پاریس بیابند.

ص: 269

1- کتابخانه به سرعت به اعتراض های بلاگرها پاسخ داد بعد از تبادل چندای میل با کتابخانه بلاگر با کاربرد قانون واسپاری وب آشنا شد و تصمیم گرفت تله اش را جمع کند

forum-2

3- توجه کنید که این تصمیم راهنمای متخصص آرشیو وب دانمارک برای انتخاب قوانین مؤدبانه محدود کننده روایات است تا مانع اضافه بار مورد درخواست سرور ها و پرونده های دعوی بالقوه شود.

4- این قفسه های Petaboxes، دارای ظرفیت بالا- هزینه پایین سخت افزار ذخیره قدرت پایین، هستند به فناوری های Capricorn اختصاص دارند. <http://www.capricorn-tech.com> [تاریخ دسترسی: 17 آوریل 2008]

3. خزش گر در حال کار

3.1. آزمایش خزش ها

قبل از اقدام به، خزش باید به منظور بررسی روپات و عکس العمل های ماشین برای پیش بینی مشکلات آزمایش خزش ها انجام شود. این، وظایف همراه با نظارت بر روپات ها در طول خزش (2/3 را ببینید) و توصیف مجموعه بعد از برداشت، عملاً مهم ترین بخش فرآیند تضمین کیفیت برای یک خزش هستند (بخش 4 را ببینید). متدولوژی، فوق برای چهار خزش فراگیر ما مورد استفاده قرار گرفت.

همان طور که قبلاً توضیح دادیم وظیفه اصلی، قبل از آغاز خزش ادغام فهرست های هسته متفاوت و بررسی آن ها بود دیگر گام مهم اجرای «آزمایش خزش» بود: برای تجزیه و کشف حجم بالای URL در فهرست هسته روپاتی راه اندازی شد اجرای خزش از آغاز تا پایان مورد نظر نیست، بلکه خزش در زمان کافی برای شناسایی اسامی دامنه غیر مرتبط یا URL های خطرناک به اجرا در می آید.

به طور مثال URL های ارسال کننده کد «خطای 404» حذف شدند به علاوه اسامی دامنه هایی که به شمار کمی میزبان عمومی (در اکثر موارد ثبت کنندگان یا مالکان تصرف کننده دامنه (1)) تغییر مسیر داده بودند نیز از خزش حذف شدند. همچنین، اسامی دامنه برای کشت دامنه (یعنی با هدف بالا بردن رتبه یک وب گاه، با استفاده از اسامی چند دامنه ای برای یک IP واحد) شناسایی شدند.

3.2. ارتباط کتابخانه ملی فرانسه و آرشیو اینترنت در طول خزش

خزش فراگیر در 11 اکتبر 2007 راه اندازی شد و در 29 نوامبر به طور کامل پایان یافت (یعنی بعد از «خزش تکه ای») در زمان خزش رابطه IA/BNF بر اساس تحلیل «گزارش یافته ها» (2) بود. این گزارش، همه دامنه های موجود در صف را با نشان دادن URL هایی که قبلاً برای هر یک از آن ها برداشت شده، مقدار بودجه صرف شده و URL های خزش شده فهرست می کند. هدف از این کار تکمیل تخصص سنتی مهندسان IA در این زمینه با دانش تیمی هسته های فرانسوی BNF است. این کار در BNF توسط یک کتابدار و یک مهندس هدایت گردید.

با افزایش بودجه به دامنه هایی که بالای 10,000 URL داشتند توجه خاصی مبذول شد. صرف نظر از اینکه داده ها مرتبط بودند یا نه این تعداد را به عنوان آستانه ای برای آن چه که باید بررسی میشد تعیین کردیم. در تعریف ما داده های غیر مرتبط اسناد با ارزش علمی اندک نیستند بلکه فایل های اضافی ای هستند که در اثر خصیصه های مربوط به وب گاه های پاتولوژی ایجاد شده اند به طور مثال، تله های روپات (به خاطر ایجاد تقویم یا نسخه های جاوا) شمار نامحدودی URL برای صفحه های وبی یکسان ایجاد کردند. وب گاه های مرور (3) نیز مشکل ساز هستند زیرا چندین دامنه محتوای یکسان را با اسامی مختلف میزبانی می کنند در اغلب موارد این دامنه ها حذف شدند.

ص: 270

1- بیش از 100000 اسامی دامنه به تنها سه وب گاه - دو ثبت نام کننده و یک متصرف دامنه تغییر مسیر داده بودند.

2- frontier report

3- Mirror websites

با استفاده از کشت دامنه وب گاه های اصلی را شناسایی و حذف کردیم. بررسی های برخط نیز برای تعیین مطابقت رشته های حروف با تقویم ها پردازش شد تا آن ها را فیلتر و از ایجاد تله های روباتی جلوگیری کند.

هدف این کنترل، هفتگی اطمینان از کسب 100 درصد کیفیت خزش نبود - به هر حال یک خزش فراگیر بود. هدف مدیریت بزرگ ترین صف ها به بهترین شکل ممکن و جلوگیری از هدر رفتن زمان و منابع زیادی توسط خزش گر بود.

3.3. «خزش تکه ای»

* «خزش تکه ای» (1)

بعد از سه هفته خزش مداوم شبانه روزی مهندسان IA تصمیم گرفتند خزش را متوقف کنند. آن ها به تحلیل QA در داده های بازبایی شده پرداختند تا دامنه هایی را مشخص کند که روبات به چهارمین سطح عمق آن ها (یعنی سه جهش نسبت به صفحه هسته) نرسیده بود این دامنه ها بار دیگر (با عنوان «خزش تکه ای») هر زمان که ممکن بود خزش شدند.

خزش تکه ای دیگری برای بازبایی فایل های ویدئویی شناخته شده در طول مدت برداشت راه اندازی شد اما بنا به دلایل فنی بارگذاری نشد (مشکلات فایل های میزبانی شده در دامنه های مختلف و جز آن).

4. پیامدهای خزش

تحلیل های زیادی درباره مجموعه برداشت شده صورت گرفت مهندسان، IA که برای کمک به ما در زمان نصب قفسه ها به پاریس آمده بودند بسیار کارآمد بودند. نخستین هدف این تحلیل ها کنترل کیفیت داده های دریافتی بود همچنین مایل بودیم مجموعه ها را در مقیاس وسیع شناسایی کنیم: طیف اسناد موجود در قفسه ها در عمل چه بود شکل و عمق وب گاه های برداشت شده چگونه بود و جز آن. از این رو، هدف، تعیین کمیت و کیفیت مجموعه سال 2007 ما بود. سرانجام تحلیل نتایج خزش فراگیر 2007، و مقایسه این نتایج با خزش های قبلی (به ویژه خزش سال 2006) برای تعیین تأثیر تنظیمات جدید خزش و تصمیم گیری درباره این که آیا این خزش ها در جهت وظیفه واسپاری قانونی ما قرار دارند، ضرورت داشت.

1.4. اشکال اصلی

عکس

با استفاده از کشت دامنه، وبگاههای اصلی را شناسایی و حذف کردیم. بررسی‌های برخط نیز برای تعیین مطابقت رشته‌های حروف با تقویم‌ها، پردازش شد تا آنها را فیلتر و از ایجاد تله‌های روباتی جلوگیری کند. هدف این کنترل هفتگی، اطمینان از کسب ۱۰۰ درصد کیفیت خزش نبود - به هر حال یک خزش فراگیر بود. هدف، مدیریت بزرگترین صف‌ها به بهترین شکل ممکن و جلوگیری از هدر رفتن زمان و منابع زیادی توسط خزش گر بود.

۳.۳. «خزش تکه‌ای»^۱

بعد از سه هفته خزش مداوم شبانه‌روزی، مهندسان IA تصمیم گرفتند خزش را متوقف کنند. آنها، به تحلیل QA در داده‌های بازبایی شده پرداختند، تا دامنه‌هایی را مشخص کند که روبات، به چهارمین سطح عمق آنها (یعنی، سه جهش نسبت به صفحه هسته) نرسیده بود. این دامنه‌ها، بار دیگر (با عنوان «خزش تکه‌ای»)، هر زمان که ممکن بود خزش شدند.

خزش تکه‌ای دیگری برای بازبایی فایل‌های ویدئویی شناخته شده در طول مدت برداشت راه‌اندازی شد، اما بنا به دلایل فنی بازگذاری نشد (مشکلات فایل‌های میزبانی شده در دامنه‌های مختلف و جز آن).

۴. پیامدهای خزش

تحلیل‌های زیادی در باره مجموعه برداشت شده صورت گرفت. مهندسان IA، که برای کمک به ما در زمان نصب قفسه‌ها به پاریس آمده بودند، بسیار کارآمد بودند. نخستین هدف این تحلیل‌ها، کنترل کیفیت داده‌های دریافتی بود. همچنین، مایل بودیم مجموعه‌ها را در مقیاس وسیع شناسایی کنیم: طیف اسناد موجود در قفسه‌ها، در عمل، چه بود، شکل و عمق وبگاههای برداشت شده چگونه بود، و جز آن. از این رو، هدف، تعیین کمیت و کیفیت مجموعه سال ۲۰۰۷ ما بود. سرانجام، تحلیل نتایج خزش فراگیر ۲۰۰۷، و مقایسه این نتایج با خزش‌های قبلی (به‌ویژه خزش سال ۲۰۰۶)، برای تعیین تأثیر تنظیمات جدید خزش و تصمیم‌گیری در باره اینکه آیا این خزش‌ها در جهت وظیفه واسپاری قانونی ما قرار دارند، ضرورت داشت.

۱.۴. اشکال اصلی

تعداد	۲۰۰۶	۲۰۰۷
URLها	۲۷۱۶۹۷۴۵۶	۳۳۷۳۲۲۲۰۰
میزبان‌ها	۲۹۲۸۳۶۴	۱۵۸۹۴۵۸
دامنه‌ها	۳۸۲۵۴۰	۱۰۶۲۳۱۷
(از دامنه های .fr)	۱۳۱۱۳۶	۷۹۱۹۴۰

1. patch-crawl

۳۱۸	۷۱۰	URL های هر دامنه
۲۱۲	۹۳	یو.آر.ال های هر میزبان
91745	73073	فایل های ARC یکتا
8,8	7,2	اندازه فشرده داده های یکتا (در Tb)

شکل ۶. اشکال اصلی خزش فراگیر fr ۲۰۰۷

به خاطر افزایش چشمگیر فهرست هسته، رشد تعداد دامنه های برداشت شده پیش بینی می شد. همکاری با آفنیک (انجمن نامگذاری اینترنت فرانسه برای ثبت fr) در سال ۲۰۰۶، به کتابخانه اجازه داد دامنه های fr گردآوری شده را در شش مرحله کشف و برداشت کند. به عبارت دیگر، در سال ۲۰۰۷ نسبت به ۲۰۰۶، با وجود افزایش یو.آر.ال های برداشت شده توسط IA، میزبان های کمتری خزش شدند. این امر، احتمال دارد به خاطر رویکرد در ازای هر دامنه، مربوط به آخرین خزش فراگیر ما باشد. شمار یو.آر.ال ها در ازای هر دامنه یا در ازای هر میزبان راهکاری ساده برای ارزشیابی «عمق متوسط» یک خزش است. این شکل، تفاوت های معنی دار فراوانی میان وبگاهها را پنهان می کند: شکل در بخش ۶.۴ با جزئیات بیشتر بررسی شده است.

۶.۴. توزیع هر سرآیند^۱

گزارش سرآیند در باره خزش فراگیر سال ۲۰۰۷، حدود ۱۶۰۴ نوع مختلف را نشان می دهد. تعجب آور نبود که یک سرآیند text/html، تنها دو سوم از فایل های برداشت شده را نمایش می دهد. به علاوه، ۹۷ درصد URL های بارگذاری شده یکی از پنج نوع سرآیند پر استفاده: HTML، JPEG، GIF، PNG و PDF، را دارند. اگر فردی سرآیندهای اسناد برداشت شده طی خزش فراگیر ۲۰۰۷ را نگاه کند، تصور می کند وب فرانسه ۲۰۰۷ را بیشتر شامل متن و تصویر است. با وجود این، باید در باره شکل ها بسیار محتاط باشیم. لازم است محدودیت های فنی رویت ها را نیز در نظر بگیریم؛ زیرا حتی اگر عملکرد خزشگر به طور مداوم پیشرفت کند، قادر به تجزیه و گردآوری قالب های مختلف فایل که در وب می باید، نیست. قالب های فایلی پیچیده، بدون نمایش بوده یا به راحتی در مجموعه غایب هستند.

دلیل مناسب دیگر برای رعایت احتیاط این است که اطلاعات سرآیند مورد نیاز محاسبه، همانی است که توسط سرور فرستاده می شود. در واقع، این اطلاعات قابل اعتماد نیست. گاهی، حتی سرور سرآیندی می فرستد که وجود ندارد (با کمال تعجب، «نرم افزار / X-چیزی»^۲ را در مجموعه خود یافتیم). به جز ۱۶۰۴ سرآیند گوناگون، ۱۴۰۰ سرآیند با کمتر از ۵۰۰ فایل مربوط هستند- می توان چنین استنباط کرد

1. MIME type
2. application/x-something

شکل 6. اشکال اصلی خزش فراگیر 2007 fr

به خاطر افزایش چشم گیر فهرست، هسته رشد تعداد دامنه های برداشت شده پیش بینی می شد. همکاری با آفنیک (انجمن نام گذاری اینترنت فرانسه برای ثبت fr) در سال 2006، به کتابخانه اجازه داد دامنه های fr گردآوری شده را در شش مرحله کشف و برداشت کند.

به عبارت دیگر در سال 2007 نسبت به 2006 با وجود افزایش یو. آر.ال های برداشت شده توسط IA میزبان های کمتری خزش شدند این، امر احتمال دارد به خاطر رویکرد در ازای هر دامنه مربوط به آخرین خزش فراگیر ما باشد.

شمار. یو. آر.ال ها در ازای هر دامنه یا در ازای هر میزبان راه کاری ساده برای ارزشیابی «عمق متوسط» یک خزش است این، شکل تفاوت های معنی دار فراوانی میان وب گاه ها را پنهان می کند: شکل در بخش 4. 6 با جزئیات بیشتر بررسی شده است.

2.4. توزیع هر سرآیند

*توزیع هر سرآیند (1)

گزارش سرآیند درباره خزش فراگیر سال 2007، حدود 1604 نوع مختلف را نشان می دهد. تعجب آور نبود که یک سرآیند text/html تنها دو سوم از فایل های برداشت شده را نمایش می دهد. به علاوه، 97 درصد URL های بارگذاری شده یکی از پنج نوع سرآیند پر استفاده HTML، JPEG، GIF، PNG و PDF را دارند اگر فردی سرآیندهای اسناد برداشت شده طی خزش فراگیر 2007 را نگاه کند، تصور می کند وب فرانسه 2007 را بیشتر شامل متن و تصویر است.

با وجود، این باید درباره شکل ها بسیار محتاط باشیم. لازم است محدودیت های فنی روایات ها را نیز در نظر بگیریم؛ زیرا حتی اگر عملکرد خزش گر به طور مداوم پیشرفت کند قادر به تجزیه و گردآوری قالب های مختلف فایل که در وب می یابد نیست قالب های فایلی، پیچیده بدون نمایش بوده یا به راحتی در مجموعه غایب هستند.

دلیل مناسب دیگر برای رعایت احتیاط این است که اطلاعات سرآیند مورد نیاز محاسبه، همانی است که توسط سرور فرستاده می شود در واقع این اطلاعات قابل اعتماد نیست گاهی، حتی سرور سرآیندی می فرستد که وجود ندارد (با کمال تعجب «نرم افزار/ X-چیزی» (2) را در مجموعه خود یافتیم). به جز 1604 سرآیند گوناگون 1400، سرآیند با کمتر از 500 فایل مربوط هستند- می توان چنین استنباط کرد

ص: 272

MIME type -1

application/x-something -2

که سرآیند ها به طرز نامناسبی تعیین شده اند.

به نظر می رسد اهمیت این مسئله سال به سال بیشتر شود. یو.آر.ال های خزش فراگیر 2004، دارای 554 سرآیند گوناگون هستند؛ این رقم به 1024 در سال 2006 و به 1604 در 2007 تغییر یافت.

عکس

قانون واسپاری وب فرانسه... ۲۷۳

که سرآیندها به طرز نامناسبی تعیین شده اند.

به نظر می رسد اهمیت این مسئله سال به سال بیشتر شود. یو.آر.ال های خزش فراگیر ۲۰۰۴، دارای ۵۵۴ سرآیند گوناگون هستند؛ این رقم به ۱۰۲۴ در سال ۲۰۰۶ و به ۱۶۰۴ در ۲۰۰۷ تغییر یافت.

درصد	یو.آر.الها	MIME-type
67.96	229257942	متن / html
19.04	64222287	تصویر / jpeg
7.52	25376262	تصویر / gif
1.17	3955885	تصویر / png
1.17	3955463	نرم افزار / pdf
۰٫۶۷	۲۲۵۶۷۵۹	متن / ساده
۰٫۴۷	۱۵۹۴۳۴۲	برنامه / فلش - shockwave-x
۰٫۴۲	1432809	متن / css
۰٫۴۲	۱۴۱۵۲۳۰	برنامه / javascript-x
0.32	1083991	برنامه / XML
0.82	2771213	سایر

شکل ۷. ده رتبه نخست سرآیند خزش فراگیر ۲۰۰۷

اما اگر نمی توانیم به سرآیند یک فایل منفرد اطمینان کاملی داشته باشیم، توزیع فراگیر تعیین شده برای صدها میلیون سند به احتمال زیاد قابل اعتماد است. تحول سرآیند می تواند به عنوان راهی برای تحلیل تغییرات و تمایلات در مقیاسی وسیع، دیده شود. به طور مثال می توان، از سال ۲۰۰۴ تا ۲۰۰۷، کاهش در استفاده از قالب GIF را، به نفع JPEG و فورمت باز PNG^۱ را مشاهده کرد (میزان تصاویر GIF در عرض این چهار سال تقریباً نصف شده است). بر همین اساس، میزان اسناد XML پنج برابر شده است. حتی وقتی به حجم اسناد پرداخت شده نگاه می کنیم رشد فایل های XML در وب و وضوح بیشتری می یابد: از ۸۷۰۰۰ در سال ۲۰۰۴ به یک میلیون در ۲۰۰۷. شاید این رشد، تا اندازه ای، به خاطر کاربرد فزاینده آر.اس.اس خوانها^۲ باشد (سرآیند صحیح برای یک فایل آر.اس.اس، «application/rss» است، اما غالباً «application/xml» یا حتی «text/xml» به جای آن به کار می روند).

۱. توجه کنید که این اشکال گاهی به چند سرآیند دسته بندی می شوند: به طور مثال، تعدادی اسناد JPEG معین با افزایش شماری اسناد دارای سرآیند «image/jpeg»، «Image/jpeg» یا «image/JPEG».

۲. توجه کنید که نرخ های مشابهی برای دامنه au مشاهده شده اند (استرالیا، ۲۰۰۴ تا ۲۰۰۵): درصد تصاویر GIF نصف شده است (از ۱۰ تا ۵ درصد)، با وجود این، تصاویر png در استرالیا (۵۶ درصد) نسبت به فرانسه (۱،۱۷ درصد) کمتر مورد استفاده قرار گرفته اند [۱۸].

3. RSS feeds

اما اگر نمی توانیم به سرآیند یک فایل منفرد اطمینان کاملی داشته باشیم توزیع فراگیر تعیین شده برای صد ها میلیون سند به احتمال زیاد قابل اعتماد است. تحول سرآیند می تواند به عنوان راهی برای تحلیل تغییرات و تمایلات در مقیاسی وسیع دیده شود. به طور مثال می توان از سال 2004 تا 2007 کاهش در استفاده از قالب GIF را به نفع JPEG و فورمت باز PNG (2) را مشاهده کرد (میزان تصاویر GIF در عرض این چهار سال تقریباً نصف شده است). بر همین اساس میزان اسناد XML پنج برابر شده است. حتی وقتی به حجم اسناد برداشت شده نگاه می کنیم رشد فایل های XML در وب وضوح بیش تری می یابد: از 88/000 در سال 2004 به یک میلیون در 2007 شاید این رشد تا، اندازه ای به خاطر کاربرد فزاینده آر.اس. اس خوان ها (3) باشد (سرآیند صحیح برای یک فایل آر.اس. اس، «application/rss» است، اما غالباً «application/xml» یا حتی «text/xml» به جای آن به کار می روند).

ص: 273

1- توجه کنید که این اشکال گاهی به چند سرآیند دسته بندی می شوند به طور مثال تعدادی اسناد JPEG معین با افزایش شماری اسناد دارای سرآیند «image/jpeg»، «Image/jpeg» یا «image/JPEG».

2- توجه کنید که نرخ های مشابهی برای دامنه au. مشاهده شده اند (استرالیا 2004 تا 2005): درصد تصاویر GIF نصف شده است (از 10 تا 5 درصد). با وجود، این تصاویر png در استرالیا (0,56 درصد) نسبت به فرانسه (1,17 درصد) کمتر مورد استفاده قرار گرفته اند

[18]

RSS feeds -3

قالب فایلی دیگر رو به افزایش در مجموعه های خزش های فراگیر ، نرم افزار فلش شاک ویو (1) است. این رشد می تواند دو علت داشته باشد: محبوبیت در حال افزایش قالب فلش در وب توانایی بیش تر خزنده هریتریکس برای برداشت این نوع فایل

عکس

۲۷۴ مدیریت منابع اطلاعاتی وب

قالب فایلی دیگر رو به افزایش در مجموعه های خزش های فراگیر ، نرم افزار فلش شاک ویو^۱ است. این رشد می تواند دو علت داشته باشد: محبوبیت در حال افزایش قالب فلش در وب، توانایی بیشتر خزنده هریتریکس برای برداشت این نوع فایل.

MIME Type evolution	2004	2005	2006	2007
text/html	68.11	67.22	70.15	67.96
image/jpeg	14.04	15.79	15.13	19.04
image/gif	12.70	11.09	8.05	7.52
application/pdf	1.36	1.39	1.19	1.17
image/png	0.79	0.73	0.87	1.17
text/plain	1.0833	1.19	1.01	0.67
application/x-shockwave-flash	0.2488	0.34	0.35	0.47
application/xml	0.07	0.16	0.50	0.32

جدول ۸. تحول در چند سرآیند از ۲۰۰۴ تا ۲۰۰۷

از منظر حفاظت بلندمدت، این اطلاعات بسیار ارزشمند است و قالبی را که - در مقیاس ملی و نیز بین المللی، باید تلاش های خود را بر آن متمرکزسازیم، ارائه می کند. کاربرد رو به افزایش قالب های باز، مانند PNG یا XML، اخبار خوبی از این نقطه نظر است.

۳. ۴. فایل های ویدئویی

افزایش چهار قالب ویدئویی پر استفاده (Quicktime, Windows media video, Flash video و MPEG video) منجر به ایجاد حدود ۴۰۰۰۰ فایل ویدئویی برداشت شده در ۲۰۰۴ گردید، در مقابل ۱۲۰۰۰۰ فایل در چهار سال بعد (یعنی ۰/۰۴ درصد مجموعه). ما متوجه کاهش قالب ویدئویی MPEG در مقابل فلش شده ایم. در ۲۰۰۶ خزشگر هریتریکس فقط صد فایل ویدئویی فلش را برداشت کرده بود، یک سال بعد، نسخه مجوزدار خزشگر ما، برای برداشت محتوای میزبانی شده در سکوها ی پخش ویدئویی به اجرا درآمد: هریتریکس سی هزار سند را گرد آورد. اراده ما برای تمرکز بر برداشت فایل های ویدئویی به هدفش رسیده بود: اگرچه اغلب آرشیوهای ما «ضعف هایی»^۲ در ویدئو داشته است، در ۲۰۰۷، نمونه ویدئویی حجیم تری را به دست آوردیم.

1. application/x-shockwave-flash
2. holes

از منظر حفاظت بلندمدت این اطلاعات بسیار ارزشمند است و قالبی را که - در مقیاس ملی و نیز بین المللی باید تلاش های خود را بر آن متمرکز سازیم ارائه می کند کاربرد رو به افزایش قالب های باز مانند PNG یا XML ، اخبار خوبی از این نقطه نظر است

3.4. فایل های ویدئویی

افزایش چهار قالب ویدئویی پر استفاده (MPEG video و Flash video Quicktime Windows media video) منجر به ایجاد حدود 40000 فایل ویدئویی برداشت شده در 2004 ، گردید در مقابل 120000 فایل در چهار سال بعد (یعنی 0/04 درصد مجموعه) ما متوجه کاهش قالب ویدئویی MPEG در مقابل فلش شده ایم. در 2006 خزش گر هریتریکس فقط صد فایل ویدئویی فلش را برداشت کرده بود، یک سال بعد، نسخه مجوزدار خزش گر ما، برای برداشت محتوای میزبانی شده در سکو های پخش ویدئویی به اجرا درآمد هریتریکس سی هزار سند را گرد آورد اراده ما برای تمرکز بر برداشت فایل های ویدئویی به هدفش رسیده بود اگر چه اغلب آرشیوهای ما «ضعف هایی» (2) در ویدئو داشته است، در 2007، نمونه ویدئویی حجیم تری را به دست آوردیم

ص: 274

application/x-shockwave-flash -1

holes -2

قانون واسپاری وب فرانسه... ۲۷۵

MIME-type	2004	2005	2006	2007
ویدئو /x-ms-wmv	4 408	7 705	33 936	39 218
ویدئو /quicktime	22 020	26 687	39 073	36 294
درخواست /x-flv	0	0	104	31 556
ویدئو /mpeg	11 408	17 304	28 413	14 992
جمع	39 840	53 701	103 532	124 067

شکل ۹. تحول سرآیند فایل‌های ویدئویی از ۲۰۰۴ تا ۲۰۰۷

۴.۴. توزیع هر TLD

دیگر کشف سه چهارم اسناد خزش شده متعلق به دامنه سطح بالای fr تعجب‌آور نبود (می‌توان به این رقم، دامنه re را که توسط آفنیق مدیریت می‌شود، افزود). با وجود این، شکل ۱۰ تأیید می‌کند که چنین چیزی با تنظیماتی که ما به کار گرفتیم، از فهرست هسته a.fr آغاز و تا مرزهای آن ادامه دادیم، ممکن است. دو نوع دامنه سطح بالا (تی.ال.دی.) ارائه شده در مجموعه نیز دامنه‌های سطح بالای کدهای عمومی و ملی^۱ هستند. اغلب سایت‌های میزبانی شده تحت دامنه‌های سطح بالای عمومی (.org، .net، .com، .info) به احتمال، توسط مدیران وب فرانسوی تولید شده‌اند. کد کشوری مربوط به رتبه‌بندی دامنه‌های سطح بالا در شکل ۱۰ یا متعلق به کشورهای فرانسوی زبان (بلژیک و سوئیس) است یا به همسایگانی که (آلمان و انگلیس) فرانسه با آنها ارتباط تجاری عمده دارد مربوط می‌شود: این پدیده برای اسپانیا نیز مورد توجه قرار گرفت [۵]. دامنه eu. ویژه اروپاست.

TLD	تعداد یوآرالها	درصد
fr	259 869 452	77.12
com	59 843 624	17.76
net	4 951 932	1.47
org	3 171 196	0.94
de	2 808 359	0.83
eu	993 456	0.29
info	900 544	0.27
be	660834	0.20
ch	461 021	0.14
uk	434 315	0.13
re	381 746	0.11
Other TLDs	2 471 064	0.73

جدول ۱۰. تعداد یوآرال‌های هر TLD، خزش‌گر فراگیر ۲۰۰۷

این ارقام، کاملاً با داده‌های مجموعه‌های قبلی، از ۲۰۰۴ تا ۲۰۰۶ مشابه است. به سقوط بالای biz. (که

1. general and country codes

شکل ۹. تحول سرآیند فایل‌های ویدئویی از ۲۰۰۴ تا ۲۰۰۷

4.4. توزیع هر TLD

دیگر کشف سه چهارم اسناد خزش شده متعلق به دامنه سطح بالای fr تعجب‌آور نبود (می‌توان به این رقم دامنه re. را که توسط آفنیق

مدیریت می شود افزود). با وجود، این شکل 10 تأیید می کند که چنین چیزی با تنظیماتی که ما به کار گرفتیم، از فهرست هسته a.fr آغاز و تا مرزهای آن ادامه دادیم، ممکن است دو نوع دامنه سطح بالا (تی.ال.دی.) ارائه شده در مجموعه نیز دامنه های سطح بالای کدهای عمومی و ملی (1) هستند اغلب سایت های میزبانی شده تحت دامنه های سطح بالای عمومی (org .net .com. info) به احتمال، توسط مدیران وب فرانسوی تولید شده اند کد کشوری مربوط به رتبه بندی دامنه های سطح بالا در شکل 10 یا متعلق به کشورهای فرانسوی زبان (بلژیک و سوییس) است یا به همسایگانی که (آلمان و انگلیس) فرانسه با آن ها ارتباط تجاری عمده دارد مربوط می شود: این پدیده برای اسپانیا نیز مورد توجه قرار گرفت. [5] دامنه eu. ویژه اروپاست

جدول 10. تعداد یو. آر. ال های هر، TLD، خزش گر فراگیر 2007

این ارقام، کاملاً با داده های مجموعه های قبلی از 2004 تا 2006 مشابه است. به سقوط بالای biz. (که

ص: 275

سال ها قبل در 10 رتبه نخست دامنه سطح بالا جای داشت)، و ظهور ناگهانی eu. توجه کنید.

با وجود این، اگر به تعداد دامنه های میزبانی شده در یک دامنه سطح بالای خاص نگاه کنیم توزیع های متفاوت را متوجه می شویم.

عکس

۲۷۶ مدیریت منابع اطلاعاتی وب

سالها قبل در ۱۰ رتبه نخست دامنه سطح بالا جای داشت)، و ظهور ناگهانی eu. توجه کنید.
با وجود این، اگر به تعداد دامنه های میزبانی شده در یک دامنه سطح بالای خاص نگاه کنیم،
توزیع های متفاوت را متوجه می شویم.

TLD	دامنه های ۲۰۰۵ (درصد)	دامنه های ۲۰۰۶ (درصد)	دامنه های ۲۰۰۷ (درصد)
com	44.59	42.78	17.10
fr	26.86	34.28	74.55
net	5.37	5.95	2.06
org	4.39	4.69	1.46

جدول ۱۱. درصد دامنه هر تی.ال.دی (منحصر به com, fr, net, و org)، از ۲۰۰۵ تا ۲۰۰۷.

مجموعه های ۲۰۰۵ و ۲۰۰۶، که با تنظیمات مشابه (فهرست هسته و تنظیمات خزش) شکل گرفته اند برتری بیشتر دامنه های سطح بالای عمومی را بر fr. نشان می دهند. این ارقام می تواند توسط حجم زیاد دامنه های غیرموجود در فهرست هسته، که توسط رویات لمس شده بود توصیف شود (یعنی جایی که یک یو.آر.ال. پیوند یافته به یک یو.آر.ال. درون حوزه، مورد خزش قرار گرفته است). برخی پیوندها به دامنه های com یا net. در صفحه های fr. قابل دسترس بودند، و به همین دلیل تا حدودی گردآوری شدند. به خاطر افزایش چشمگیر اندازه فهرست هسته^۱، میزان زیادی از دامنه های قابل دسترس در مجموعه ۲۰۰۵ و ۲۰۰۶ - اما نه در خزش فراگیر ۲۰۰۷ - عرضه می شوند.

۴. Robots.txt. ۵.

در سال ۲۰۰۷، پروتکل حذف رویات^۲ مانع ما در آرشوسازی ۱۵ میلیون فایل بود، یعنی، ۴.۵ درصد فایل های کشف شده - و به وضوح از آرشوسازی تمام اسنادی که رویات می توانست از آغاز تشکیل این فایل ها کشف کرده باشد جلوگیری کرد. این ارقام درکل، زمانی که فایل های robots.txt بارگذاری ۶ درصد اسناد کشف شده را متوقف کردند نسبت به خزش فراگیر ۲۰۰۶ بسیار پایین تر هستند. تفسیر این میزان مشکل است، زیرا آنها با آخرین مطالعات این پروتکل، که کاربرد در حال رشد robots.txt را شناسایی می کند، مغایرت دارند [به طور مثال منبع ۲۳].

زمان تحلیل آنچه خزش نشده بود، استفاده از سرآیند فایل ها امکان نداشت، زیرا سرور آنها را ارسال نکرده بود. با وجود این، می توانیم از پسوند فایل URL های خواسته شده استفاده کنیم. تقریباً ۴۰ درصد این فایل ها تصویر هستند^۳. ممکن است برخی مدیران وب می خواستند مانع خزش رویات در این فایل ها

۱. توجه کنید برخی دامنه های برداشت شده com یا net بین ۲۰۰۶ و ۲۰۰۷ تفاوت خیلی زیادی با یکدیگر ندارند. ۱۶۳۶۳۲ دامنه com در ۲۰۰۶ نسبت به ۱۸۱۶۲۶ در ۲۰۰۷؛ ۲۲۷۴۶ دامنه net در ۲۰۰۶ نسبت به ۲۱۸۵۳ در ۲۰۰۷.
۲. robots Exclusion Protocol
۳. jpg. (۲۶ درصد)، gif. (۸ درصد)، JGP. (۲ درصد)، png. (۲ درصد). به علاوه توجه کنید که ۱۰۴،۱۰۹ فایل (۱ درصد) از فایل های CSS بودند.

جدول 11. درصد دامنه هر تی.ال.دی (منحصر به com, fr, net, و org)، از 2005 تا 2007.

مجموعه های 2005 و 2006، که با تنظیمات مشابه (فهرست هسته و تنظیمات خزش) شکل گرفته اند برتری بیش تر دامنه های سطح بالای عمومی را بر fr نشان می دهند این ارقام می تواند توسط حجم زیاد دامنه های غیر موجود در فهرست، هسته که توسط روبات لمس شده بود توصیف شود (یعنی جایی که یک یو.آر.ال. پیوند یافته به یک یو.آر.ال درون حوزه مورد خزش قرار گرفته است). برخی پیوندها به دامنه های com. یا net. در صفحه های fr قابل دسترس بودند و به همین دلیل تا حدودی گردآوری شدند به خاطر افزایش چشمگیر اندازه فهرست هسته (1) میزان زیادی از دامنه های قابل دسترس در مجموعه 2005 و 2006 - اما نه در خزش فراگیر 2007 عرضه می شوند.

Robots.txt.5.4

در سال 2007، پروتکل حذف روبات ها (2) مانع ما در آرشیو سازی 15 میلیون فایل بود، یعنی، 4,5 درصد فایل های کشف شده - و به وضوح از آرشیو سازی تمام اسنادی که روبات می توانست از آغاز تشکیل این فایل ها کشف کرده باشد جلوگیری کرد این ارقام در کل زمانی که فایل های robots.txt بارگذاری 6 درصد اسناد کشف شده را متوقف کردند نسبت به خزش فراگیر 2006 بسیار پایین تر هستند. تفسیر این میزان مشکل است زیرا آن ها با آخرین مطالعات این پروتکل که کاربرد در حال رشد robots.txt را شناسایی می کند مغایرت دارند [به طور مثال منبع 23].

زمان تحلیل آن چه خزش نشده بود استفاده از سرآیند فایل ها امکان نداشت زیرا سرور آن ها را ارسال نکرده بود. با وجود این می توانیم از پسوند فایل URL های خواسته شده استفاده کنیم. تقریباً 40 درصد این فایل ها تصویر هستند (3) ممکن است برخی مدیران وب می خواستند مانع خزش روبات در این فایل ها

ص: 276

1- توجه کنید برخی دامنه های برداشت شده com. یا net بین 2006 و 2007 تفاوت خیلی زیادی با یکدیگر ندارند. 163632 دامنه com. در 2006 نسبت به 181626 در 2007؛ 22746 دامنه net. در 2006 نسبت به 21853 در 2007

2- robots Exclusion Protocol

3- (26 jpg درصد) 8، (gif درصد) 2 (JGP درصد) 2 (png درصد) به علاوه توجه کنید که 104,109 فایل (1 درصد) از فایل های CSS بودند.

شوند، زیرا آن‌ها توسط موتورهای جست‌وجو نمایه‌سازی نشده بودند این فرض با کمک شیوه‌ای که روبات‌های ما این فایل‌ها را کشف کردند مورد تأیید قرار گرفته است: نیمی از آن‌ها (7/378/578)، هنگام پیگیری یک پیوند جاسازی شده یافت شدند.

بنابراین، تمکین از robots.txt مانع برداشت بسیار مرتبط داده‌هایی می‌شود که برای روبات‌های موتورهای جست‌وجو، غیر ضروری اما برای روبات‌آرشیوسازی ما بسیار مفید است. در برخی مواقع، پروتکل حذف روبات‌ها (آر.ای.پی) از دسترسی ما به کل یک سایت جلوگیری می‌کرد بیش از URL 150/000 که به عنوان هسته‌هایی که به وسیله robots.txt محافظت می‌شدند به کار رفتند.

4.6 عمق خزش

برای تعیین عمق خزش، ممکن است تعداد یو. آر. ال‌ها برای هر دامنه .fr را محاسبه کنیم.

عکس

شوند، زیرا آنها توسط، موتورهای جست‌وجو نمایه‌سازی نشده بودند. این فرض با کمک شیوه‌ای که روبات‌های ما این فایل‌ها را کشف کردند مورد تأیید قرار گرفته است: نیمی از آنها (۷/۳۷/۵۷۸)، هنگام پیگیری یک پیوند جاسازی شده، یافت شدند.

بنابراین، تمکین از robots.txt مانع برداشت بسیار مرتبط داده‌هایی می‌شود که برای روبات‌های موتورهای جست‌وجو، غیرضروری، اما برای روبات آرشیوسازی ما بسیار مفید است. در برخی مواقع، پروتکل حذف روبات‌ها (آر.ای.پی) از دسترسی ما به کل یک سایت جلوگیری می‌کرد: بیش از ۱۵۰/۰۰۰ URL که به‌عنوان هسته‌هایی که به وسیله robots.txt محافظت می‌شدند به کار رفتند.

۶.۴ عمق خزش

برای تعیین عمق خزش، ممکن است تعداد یو.آر.ال‌ها برای هر دامنه .fr را محاسبه کنیم.

تعداد دامنه‌ها	تعداد یو.آر.ال‌ها
498777	10<
146356	10-100
103370	100-1000
43101	1000-10000
334	10000>

جدول ۱۲. تعداد یو.آر.ال‌ها برای هر دامنه .fr، خزش فراگیر ۲۰۰۷

تقریباً ۵۰ درصد دامنه‌های برداشت شده شامل ۱۰ یو.آر.ال یا کمتر هستند. این امر می‌تواند به سبب عدم دسترسی به سرور دائمی در طول خزش باشد. با وجود این، ما چندین آزمایش «دستی» انجام دادیم، یعنی برای کشف آنچه به‌صورت برخط قابل دسترس بود، در یک دو جین نام دامنه با آستانه ۱۰ یو.آر.ال. کلیک کردیم. این آزمایش‌ها نشان دادند که این وبگاه‌ها خالی بودند (مالکان آنها را خریداری کردند تا مطمئن شوند که آنها استفاده نخواهند شد، اما خودشان هم از آنها استفاده نمی‌کنند) یا اینکه آنها برای کشت پیوند استفاده شدند.

این ارقام با ارقام متعلق به خزش فراگیر قبلی بسیار تفاوت داشتند.

تعداد دامنه‌ها	تعداد یو.آر.ال‌ها
38 439	10<
32 258	10-100
41 352	100-1000
15 159	1000-10000
3 928	10000>

شکل ۱۳: تعداد یو.آر.ال‌ها برای هر دامنه .fr، خزش فراگیر ۲۰۰۶

جدول 12. تعداد یو.آر.ال‌ها برای هر دامنه .fr، خزش فراگیر 2007

تقریباً 50 درصد دامنه‌های برداشت شده شامل 10 یو.آر.ال یا کم تر هستند. این امر می‌تواند به سبب عدم دسترسی به سرور دائمی در طول خزش باشد با وجود این ما چندین آزمایش «دستی» انجام دادیم یعنی برای کشف آن چه به صورت برخط قابل دسترس بود در یک دو جین نام دامنه با آستانه 10 یو.آر.ال کلیک کردیم این آزمایش‌ها نشان دادند که این وبگاه‌ها خالی بودند (مالکان آن‌ها را خریداری کردند تا مطمئن شوند که آن‌ها استفاده نخواهند شد، اما خودشان هم از آن‌ها استفاده نمی‌کنند) یا این که آن‌ها برای کشت پیوند استفاده شدند.

این ارقام با ارقام متعلق به خزش فراگیر قبلی بسیار تفاوت داشتند.

شکل 13: تعداد یو.آر.ال ها برای هر دامنه .fr. خزش فراگیر 2006

ص: 277

سه تفاوت عمده میان سال های 2006 و 2007 درباره عمق دامنه های .fr عبارت اند از:

- دامنه های بسیار زیاد تقریباً خالی در 2007 که می توانست ناشی از افزایش اندازه فهرست هسته باشد.

- برخی دامنه ها در جایی که بیش از 10,000 یو آر ال آرشیو شده بود میان دو خزش به عدد 10 تقسیم شدند این امر پیامد محدود کردن «بودجه» برای 10,000 یو.آر.ال بود؛ همچنین شاید به خاطر رویکرد «به ازای هر دامنه»: دامنه های چندین میزبان باز نمون نمی شوند.

- از سوی دیگر این رویکرد خزش عمیق تر تعداد بیش تر دامنه های «کوچک» یا «متوسط»، بین 100 و 10,000 یو.آر.ال را اجازه می داد.

ممکن است این ارقام تحت تأثیر تله های رویات، باشند اما در حال حاضر نمی دانیم با کدام پسوند. سایر کنترل کیفیت ها، بیش از همه کنترل بصری وب گاه های شخصی برای تخمین این ارقام ضرورت دارد. با وجود این ارقام از دید قانون واسپاری فرانسه و اهداف اولیه این خزش، رضایت بخش به نظر رسد: با تضمین اینکه همه وب گاه های .fr موجود در مجموعه وظیفه سنگین برداشت وب گاه های کوچک و متوسط شاید به بهای وب گاه های بزرگ تر را داشته باشند

7.4. وب گاه های بزرگ

برای پالایش تحلیل عمق وب گاه ممکن است بر بزرگترین وب گاه ها تمرکز کنیم. هدف، بررسی علت چرایی بیش تر خزش شدن این وب گاه ها نسبت به دیگران است - آیا به این دلیل است که آن ها نسبت به دیگران بیش تر برخط هستند؟

اطلاعات مختلفی از نمایه خزش (سی.دی.ایکس) (1) اخذ شده بود: فهرست 50 دامنه بزرگ (مربوط به سال 2005 تا 2007)؛ و فهرست 1000 دامنه بزرگ در سال های 2006 و 2007.

7.4.1. دامنه ها

از این گذشته میان مجموعه های مختلف تناقص های زیادی وجود دارد نخستین تناقض، اندازه دامنه های بزرگ است تنها سه دامنه دارای بیش از 1,000,000 یو آر ال در مجموعه در 2007 نسبت به 25 یو.آر.ال در مجموعه 2006 - در بخش قبلی ارقام متفاوتی را دیدیم و می توانیم با همان دلایل این ارقام را توصیف کنیم، هم چنین محتوای 50 سایت بزرگ نخست فهرست بسیار متفاوت است: تنها 20 درصد 50 دامنه بزرگ نخست مربوط به خزش گرهای فراگیر سال 2007 در وضعیتی معادل آن در 2006 ارائه شده اند.

بزرگ ترین وبگاه 2006 (free.fr) که دارای بیش از هفت میلیون یو.آر.ال بود، در سال 2007 فقط 40000 یو.آر.ال «وزن» دارد!

ص: 278

قانون واسپاری وب فرانسه... ۲۷۹

تعداد یو.آرال.ها	نام دامنه ۲۰۰۶	تعداد یو.آرال.ها	نام دامنه ۲۰۰۷
7 405 987	free.fr	3 984 821	asso.fr
5 194 030	amiz.fr	1 760 957	com.fr
4 036 224	asso.fr	1 270 244	tm.fr
3 482 657	lrencontre.fr	534 599	gouv.fr
2 547 231	sportblog.fr	495 753	cci.fr
2 360 960	gouv.fr	408 881	co.uk
2 113 314	promovacances.fr	179 256	nom.fr
1 895 302	football.fr	144 207	presse.fr
1 885 549	mbpro.fr	108 222	dailymotion.com
1 856 720	com.fr	102 081	notaries.fr

شکل ۱۴. ده دامنه بزرگ اول، از ۲۰۰۵ تا ۲۰۰۷

به علت انجام تنظیمات خاص در مجموعه ۲۰۰۷، بالاترین رتبه دامنه‌ها تقریباً مربوط به دامنه‌های سطح دوم است. از سوی دیگر، در خزش‌های ۲۰۰۵ و ۲۰۰۶، دو نوع وبگاه کشف شد: سکوهایی با میزبانی و بنوشت‌ها و وبگاه‌های شخصی (free.fr, sportblog.fr) که توسط رویکرد به ازای هر میزبان پشتیبانی شده‌اند، و وبگاه‌های تجاری با آگهی در صفحه‌های متعدد (به‌طور مثال promovacances.fr, آرژانس مسافرتی برخط). خزش فراگیر ۲۰۰۵ نیز چند وبگاه دانشگاهی را نشان می‌دهد (۱۶ وبگاه دانشگاهی در ۵۰ دامنه بزرگ اول)، مانند jussien.fr یا cnrs.fr. در سال‌های بعد، این وبگاه‌ها تقریباً در فهرست ۵۰ وبگاه بزرگ اول وارد شدند.

حضور گسترده وبگاه‌های تجاری توصیف‌کننده تعداد دامنه‌های سطح بالای عمومی در ۵۰ دامنه بزرگ اول هستند: حتی برای مجموعه ۲۰۰۷، تنها ۳۲ درصد در fr هستند. به این ترتیب، توصیف مجموعه منعکس‌کننده وب فرانسه در سال ۲۰۰۷ چنین است: عمدتاً، فضایی برای تجارت، خدمات و روابط اجتماعی.

۴. ۷. ۲. دامنه‌های سطح دوم

	2006	2007	Evolution
asso.fr	4 036 224	3 984 821	↙
com.fr	1 856 720	1 760 957	↙
tm.fr	1 150 555	1 270 244	↗
gouv.fr	2 360 960	534 599	↙

شکل ۱۵. تحول دامنه‌های سطح دوم از ۲۰۰۶ تا ۲۰۰۷

شکل ۱۴. ده دامنه بزرگ اول، از ۲۰۰۵ تا ۲۰۰۷

به علت انجام تنظیمات خاص در مجموعه ۲۰۰۷، بالاترین رتبه دامنه‌ها تقریباً مربوط به دامنه‌های سطح دوم است از سوی دیگر در خزش‌های ۲۰۰۵ و ۲۰۰۶، دو نوع وبگاه کشف شد سکوهایی با میزبانی و بنوشت‌ها و وبگاه‌های شخصی (free.fr, sportblog.ir) که توسط رویکرد به ازای هر میزبان پشتیبانی شده‌اند و وبگاه‌های تجاری با آگهی در صفحه‌های متعدد (به‌طور مثال

fr. promovacances، آژانس مسافرتی برخط) خزش فراگیر 2005 نیز چند وب گاه دانشگاهی را نشان می دهد (16 وب گاه دانشگاهی در 50 دامنه بزرگ اول)، مانند jussien.fr یا cnrs.fr. در سال های بعد این وب گاه ها تقریباً در فهرست 50 وب گاه بزرگ اول وارد شدند.

حضور گسترده وب گاه های تجاری توصیف کننده تعداد دامنه های سطح بالای عمومی در 50 دامنه بزرگ اول هستند: حتی برای مجموعه 2007 تنها 32 درصد در fr. هستند به این ترتیب، توصیف مجموعه منعکس کننده وب فرانسه در سال 2007 چنین است: عمدتاً فضایی برای تجارت، خدمات و روابط اجتماعی

2.7.4. دامنه های سطح دوم

شکل 15. تحول دامنه های سطح دوم از 2006 تا 2007

ص: 279

توجه خاص به دامنه های سطح دوم در طول برداشت سال 2007 به کتابخانه ملی فرانسه خزش مقادیر داده اندکی پایین تر را نسبت به مشابه آن ها در برداشت سال 2006 اجازه داد. مهم ترین استثنا دامنه سطح دوم .gou.fr است از 2/3 میلیون به 500,000 یو.آر.آل سقوط کرده است. یک راه برای توصیف این رقم کاربرد مشترک دامنه های سطح سوم و حتی چهارم در وب گاه های دولتی (به طور مثال، www.rhone.pref.gouv.fr یا www.auvergne.culture.gouv.fr) است.

تأکید بر فایل های ویدئویی مجموعه بهتری از وب گاه های پخش ویدئویی (1) به بار آورد دیلی موشن (2) 10 تا از پر رتبه ترین ها دامنه ها را وارد کرد (در سال 2006، تنها 30,000 یو.آر.آل در این دامنه برداشت شده بود). تعداد فایل های گردآوری شده در یوتیوپ دوبرابر شدند.

5. نتیجه گیری

برای برداشت دامنه ملی فقط یک راه وجود ندارد. انتخاب های فنی متفاوتی (قبل و حین خزش، و حتی بعد از حذف URL ها) مجموعه را شکل می دهند سیاست مجموعه سازی حتی در مقیاس وسیع اعمال می شود. برای اجرای یک، خزش دامنه بر طبق قانون و اهداف، آن شناسایی این انتخاب ها و تخمین پیامد های آن ها ضروری است.

به طور سنتی در رسالت قانون و اسپاری برای تصمیم گیری درباره مرتبط بودن سند به حوزه، مجموعه سه معیار به کار می رفته است باید با قالبی، خاص قابل دسترس مردم و داخل مرزهای سرزمین فرانسه موجود باشد، همه این ها باید ویژگی های جدید وب را پذیرفته باشند و باید هنگام انجام خزش فراگیر به حساب آمده باشند.

هدف برداشت وب «فرانسه» تأکید بر .fr را توصیف می کند در واقع این امر برآورده نمی شود اگر تصور کنیم 50 تا 60 درصد وب گاه های فرانسوی خارج از .fir هستند اما در حال حاضر تأکید بر .fr. یک انتخاب واقع گرایانه و اقتصادی است همه محتوای فرانسه در .fr موجود نیست، اما هر چیزی که در .fr است متعلق به فرانسه است. به علاوه این تأکید قابل انعطاف است زیرا روایات اجازه دارد از تغییر مسیرهای .fir به دیگر دامنه ها تبعیت کند. سرانجام دامنه .fr به سرعت با ساده کردن قوانین ویژه این دامنه در حال توسعه است و امید می رود به زودی بخش بزرگ تری از دامنه فرانسه را بنمایاند. انتظار می رود این خواسته با قوانین CCTLD جدید مربوط به ICANN متوقف نگردد.

همچنین تأکید بر تأکید بر .fr بسیار آسان است زیرا مطمئن هستیم به لطف موافقت نامه با آفنیک، قادر به برداشت همه جانبه این دامنه سطح بالا هستیم این راه کاری است که با دومین اصل مهم قانون واسپاری منطبق است: گردآوری کل تولید فرهنگی کشور هر آن چه به محض قرار گرفتن در دسترس عموم «کیفیت» می یابد یا «ارزش» دارد. از تعداد بسیار زیاد هسته آغاز کردن ضمانتی است برای فراموش نکردن وب گاه های با پیوند کمتر یا نامشهور

ص: 280

این اصول «غیر تبعیض آمیزانه» با ویژگی دیگر قانون واسپاری: میل به گردآوردن تمام قالب های انتشاراتی در حال ظهور ناسازگار نیستند. کتابخانه در گذشته [از قالب های مختلفی] پشتیبانی می کرد: متون، تصویر، صدا یا ویدئو. با آرشیوسازی وب دیگر رفتار متفاوت با آن ها قابل تصور نیست زیرا این قالب ها به عناصری در شبکه ای یکسان پیوند یافته اند با وجود، این در صورت ضرورت می توان به راه حل های مناسبی برای انواع مختلف رسانه دست یافت گردآوری نمایه سازی، حفاظت و دستیابی به فایل های متنی و دیداری - شنیداری یا تعاملی همواره به مسائل یکسانی ختم نمی شوند.

هنوز پرسش هایی باقی است: به طور مثال دشوار است که بگوییم، در حال حاضر، برای کسب تصویری جامع از وب فرانسه آیا رویکرد تأکید تنها بر دامنه ها بهتر از نگهداری به ازای هر میزبان به طور جداگانه است لازم است تحلیل بیشتری برای پاسخ به این سؤال انجام شود و کتابخانه ملی فرانسه، به شنیدن گزارش های بین المللی درباره این موضوع بسیار علاقه دارد سرانجام، مسئله، چگونگی تعریف یک وبگاه است زیرا این هویت هوشمند اغلب با میزبانی فنی سازگار نیست اگر وب گاه را به عنوان دامنه تعریف کنیم، رویکرد به ازای هر دامنه باید پذیرفته شود، زیرا به خزش بهتر وب گاه های کوچک یا متوسط منتهی می شود. اما اگر وبگاه را به عنوان هویتی فکری ایجاد شده توسط یک نویسنده یا یک ویراستار (یک یا چند شخص یک مؤسسه عمومی یا خصوصی) تعریف کنیم سایت و دامنه دیگر با هم سازگار نیستند. در واقع، تعداد یا حتی حجم عظیمی از وب گاه ها می توانند تحت یک نام دامنه میزبانی شوند، همچون و نوشتن های میزبانی شده در سکوها تجاری. این سایت ها - با وجود مرتبط بودن - با راه کاری که در سال 2007 استفاده کردیم باز نمون نشد «به طور مثال، صفحه های شخصی فراوانی، به کل، از free.fr در 2006 برداشت شدند اما یک سال بعد ناپدید شدند برای اجتناب از این مسئله، در خزش های فراگیر آینده به کارگیری رویکردی خاص برای بسیاری از سکوها عمومی برای میزبانی کردن برخی وب گاه ها یا وب نوشت های شخصی امکان پذیر است زیرا متوجه سکوها پخش ویدئو در 2007 بودیم.

این تصمیم ها نشان دهنده اصلاحاتی است که می توانست هدف خزش هریتریکس باشد. قابلیت های بهتر برای تجزیه و برداشت قالب های پیچیده، فایل یکی از مسائل ضروری است - از این نقطه نظر هریتریکس پیش از این خود را بسیار پیکربندی شده نشان داد. در چارچوب پروژه «خزشگر هوشمند»، تاکنون سه ویژگی، دیگر توسط IA، IIPC، کتابخانه کنگره کتابخانه، بریتانیا و کتابخانه ملی فرانسه - که هدف شان توسعه روایات هریتریکس است - ارائه شده است نخستین، آن ها اجتناب از برداشت محتوایی است که از زمان آخرین خزش تغییری نکرده است: این ویژگی کاهش تکرار به صرفه جویی در تخمین منابع و ذخیره سازی آن ها و بنابراین به خزش عمیق تر وب گاه ها منجر می شود. دومین ویژگی، دادن مجوز به روایات برای اولویت بندی URL های داخل صف است. تلفیق رویکرد تمام گزینشی (در آغاز خزش) با به کارگیری قابلیت های کنترل خزش خودکار بهینه بسیار مفید خواهد بود. سومین توسعه، شناسایی خودکار بسامد تغییر وب گاه ها نیز به روایات اجازه شناسایی سایت هایی که باید به آن ها توجه خاصی مبذول شود، می دهد. از این رو، به طور مثال این امر باعث می شود تاریخ ها و بسامدهای خزش فراگیر را انتخاب کنیم، یا خزش های کانونی را در سایت های پر تغییر اجرا کنیم

در حقیقت اگر قرار است خزش های فراگیر با عمقی متوسط هر وب گاهی را یک یا دوبار در سال برداشت کنند باید هدف از خزش های کانونی را چنین تعریف کنیم: خزش های کانونی باید از ابتدا قصد شان آرشیو وب گاه های بزرگ و عمیق باشد؛ نه وب گاه های فرانسوی fit یا وب گاه های پرتغییر - آرشیو آن ها به بهترین شکل ممکن زیرا حتی خزش های کانونی اغلب برای برداشت کامل وب گاه های عظیم و گردآوری اسناد موجود در وب پنهان کافی نیستند با وجود این باید به خاطر داشته باشیم که تهیه تصاویر تنها راه - و اقتصادی ترین راه - گردآوری حافظه دیجیتال فرانسه است و اتخاذ هر تصمیم در این موضوع در سایر راه های آرشیوسازی در راهبرد تلفیقی باید به حساب آیند.

تشکر و قدردانی

مایلیم از کریس کارپنتر (1) و تیم آرشیو، اینترنت به عنوان همکار ما در این چهار سال پایانی (از ابتدا!) قدردانی کنیم به ویژه ایگور رانیتوویک (2) که بر تمام خزش های BNF مربوط به سال های 2004 تا 2007 نظارت داشت همراه با جان لی (3)، برد تافل (4) و میخائیل ماگین (5) که کمک کردند قفسه ها را نصب کنیم و قالب های نوین مجموعه را در پاریس و سان فراسیسکو تحلیل نماییم.

همچنین سپاس بسیار از مارالینو چوونیک (6) برای دقت در ویرایش سرانجام، مراتب سپاس مان را به گیلداس ایلین (7)، رئیس واسپاری دیجیتال به خاطر توصیه ها و حمایت هایش تقدیم می کنیم

ص: 282

Kris Carpenter -1

Igor Ranitovic -2

John Lee -3

Brad Tofel -4

Michael Magin -5

Mireille Chauveinc -6

Gildas Illien -7

- Abiteboul, S., Cobena, Masanè's, J. and Sedrati, G. 2002. A First Experience in Archiving the French [1] Web. In Proceedings of the Research and advanced technology for digital libraries: 6th European conference .(Italy, 2002
- AFNIC. 2007. French Domain Name Industry report. 2007 Edition. AFNIC, Saint Quentin en Yvelines. [2]
<http://www.afnic.fr/data/actu/public/2007/afnic-frenchdomain-name-report-2007.pdf>
- Andersen, B. 2005. The DK-domain: in words and figures. Netarkivet.dk, Aarhus, Copenhagen. [3]
<http://netarchive.dk/publikationer/DFreyv-english.pdf>
- Ashenfelder, M. 2006. Web Harvesting and Streaming Media. In Proceedings of the 6th International [4] Web Archiving Workshop (Alicante, Spain). <http://www.iwaw.net/06/PDF/iwaw06-proceedings.pdf>
- Baeza-Yates, R., Castillo, C. and Lopez, V. 2005. Characteristics of the Web of Spain. In Cybermetrics, [5] 9. <http://www.catedratelefonica.upf.es/webes/2005/Characteristics Web-Spain.pdf>
- Baeza-Yates, R., Castillo, C., Marin, M. and Rodriguez, A. 2005. Crawling a country: Better Strategies [6] than BreadthFirst for Web Page Ordering. In Proceedings of the 14th international conference on World .(Wide Web (Chiba, Japan
- Baly, N. and Sauvin, F. 2006. Archiving Streaming Media on the Web, Proof of concept and Firsts [7] Results. In Proceedings of the 6th International Web Archiving Workshop (Alicante, Spain).
<http://www.iwaw.net/06/PDF/iwaw06-proceedings.pdf>
- Brin, S. and Page, L. 1998. The Anatomy of a Large-scale Hypertextual Web Search Engine. In [8] Computer Networks and ISDN Systems, 30 (1-7), 107-117. <http://www7.scu.edu.au/programme/fullpapers/1921/com1921.htm>
- Dailymotion. Dailymotion [Accessed: May 10, 2008]. Partagez vos vidéos. [9]
<http://www.dailymotion.com>
- Gomes, D and Silva, M. Characterizing a National Community Web. ACM Transactions on Internet [10] Technology (volume 5, issue 3), New York, 508-531. <http://xldb.fc.ul.pt/daniel/gomesCharacterizing.pdf>
- [Heritrix. Heritrix Home Page. <http://crawler.archive.org> [Accessed: May 22, 2008 [11]

- IIPC. International internet preservation consortium-welcome. <http://www.netpreserve.org>. [13]
.[[Accessed: May 24,2008
- Illien, G., Aubry, S., Hafri Y. and Lasfargues, F 2006. Sketching and checking quality for web archives: [14]
a first stage report from BnF. Bibliothèque nationale de France, Paris.
<http://bibnum.bnf.fr/conservation/index.html>
- Illien, G. 2006. Web archiving at BnF. In International Preservation News, Paris, BnF, 27-34. [15]
<http://www.ifla.org/VI/4/news/ipnn40.pdf>
- Kimpton, M., Braggs, M. and Ubois, J. 2006. Year by Year: From an Archive of the Internet to an [16]
.Archive on the Internet. In Web Archiving, J. Masanè s, Ed, Springer, Berlin, Heidelberg, New York
- Koerbin, P. 2005. Report on the crawl and Harvest of the Whole Australian Web Domain Undertaken [17]
during June and July 2005. National Library of Australia, Canberra. [http://
pandora.nla.gov.au/documents/domain-harvest-report-public.pdf](http://pandora.nla.gov.au/documents/domain-harvest-report-public.pdf)
- Koerbin, P. 2008. The Australian Web domain harvests: a preliminary quantitative analysis of the [18]
archive data. National Library of Australia, Canberra. <http://pandora.nla.gov.au/documents/auscrawls.pdf>
- Masanè s, J. 2002. Towards continuous Web Archiving: First results and an agenda for the future. In D- [19]
Lib Magazine, 8 (12).<http://www.dlib.org/dlib/december02/masanes/12masanes.html>
- Masanè s, J. 2006. Selection for Web Archives. In Web Archiving, J. Masanè s, Ed, Springer, Berlin, [20]
.Heidelberg, New York
- Mohr, G., Kimpton, M., Stack, M, and Ranitovic, I. 2004.Introduction to Heritrix, an archival quality [21]
Web crawler.Paper presented at the 4th International Web Archiving Workshop (Bath, United Kingdom,
2004). <http://www.iwaw.net/04/Mohr.pdf>
- Najork, M. and Wiener, J. L. 2001. Breadth-First Search Crawling Yields High-Quality Pages. In: [22]
Proceedings of the 10th international conference on World Wide Web. Elsevier Science, Hong Kong, 114-
.118
- Sun, Y. Zhuang Z., Council I. and Giles C L. 2007 Determining Bias to Search Engines from [23]
Robots.txt. In Proceedings of the IEEE/WIC/ACM International Conference on Web Intelligence. IEEE
Computer Society Washington, 149-155. <http://www.personal.psu.edu/yus115/docs/sun-robotstxtbias.pdf>

:Sun, Y. Zhuang Z. and Giles C. L..2007. A large-scale study of robots.txt. In WWW '07 [24]

ص: 284

Proceedings of the 16th international conference on World Wide Web, ACM Press, New York 1123-1124.

http://www2007.org/posters/poster_1034.pdf

YouTube. YouTube - Broadcast Yourself. <http://www.youtube.com> [Accessed: May 15, 2008 [25]]

ص: 285

مجموعه های کتابخانه ملی فرانسه (1) بخشی از میراث ملی هستند و به طور تقریبی 31 میلیون سند از همه نوع (کتاب نشریه نسخه های خطی عکس ها، نقشه ها و غیره) را شامل می شوند. چالش های جدید مجموعه با گسترش اینترنت ایجاد شده اند. کتابخانه ملی فرانسه در قالبی بین المللی رهنمون های خط مشی، گردش کارها، و ابزارها را توسعه می دهد تا قسمت های مرتبط و معرف بخش فرانسوی اینترنت را جمع آوری کند و حفاظت و دسترسی به آن ها را سازماندهی کند.

آرشیوهای وب حوزه ملی فرانسه به عنوان خدمتی جدید توسعه یافت به عنوان کاربردی جدید عرضه شد و در آوریل 2008 در دسترس عموم قرار گرفت. از آن، پس راهبردهایی بوده است و توسعه می یابد تا کتابداران را درگیر آن کند و آن را در دسترس کاربران نهایی قرار دهد

این، مقاله تجربه کتابخانه ملی فرانسه را به ویژه با تمرکز بر این چهار موضوع بررسی می کند:

- ساختمان مجموعه آرشیو وب به عنوان مجموعه ای جدید و چالش برانگیز،
 - کشف منبع: خدمات و ابزارهای دسترسی برای کاربران نهایی
 - کاربرد: اطلاعات و ارقام
 - مشارکت: راهبردهایی برای ساخت یک انجمن کتابداران و رسیدن به کاربران نهایی
- کلید واژه ها: آرشیوهای ملی آرشیو کردن وب گاه ها، ساختمان مجموعه، کشف منبع، کاربران نهایی، کاربرد، فرانسه

ص: 286

نوشته: سارا اوبری (1) | ترجمه: زهرا تهوری (2)

ساختمان مجموعه آرشیوهای وب به عنوان مجموعه ای جدید و چالش برانگیز اینترنت نقش مهمی در زندگی های روزانه ما به عهده گرفته است به عنوان مثال، مدیریت الکترونیکی یادگیری الکترونیکی تجارت الکترونیکی انتشارات آنلاین هنرهای دیجیتال بلاگ ها و فضاهای جدید عمومی بحث و گفتگو بسیاری از فعالیت ها تاحدی یا به طور کامل به سمت وب حرکت کرده و تغییر مکان داده و فعالیت های جدیدی ایجاد کرده. اند با افزایش مستمر و رو به رشد تعداد کاربران اینترنت (امروزه، حدود 35 میلیون در فرانسه) و تعداد فزاینده وب گاه های فرانسوی (فقط 1/7 میلیون دامنه .fr. به نام Top Level Domain or TLD ثبت می شود) (3)، حیاتی است که این نوع انتشارات و رسانه ارتباطی را که هنوز تاحدی برای یک کتابخانه ملی جدید است، بررسی کرد.

یک وبگاه در موارد زیر متفاوت از دیگر انواع انتشارات است:

- محدود به شکل خاصی همچون یک تکه کاغذ یا یک قطعه موسیقی روی دیسک نیست بلکه وابسته به زیرساخت شبکه ای پیچیده تری است یک وبگاه یک فایل پی.دی.اف تنها (شبهه پایان نامه های

ص: 287

1- Sara Aubry, Web Archiving Project Manager, IT Department, National Library of France, Quai François – 1

Mauriac, sara.aubry@bnf.fr

2- کارشناس سازمان اسناد و کتابخانه ملی ایران

3- Observatoire du marché des noms de domaine en France, FR Network Information Center, –

[http://www.afnic.fr/actu/observatoire [last accessed on 2010-06-15]

الکترونیکی) یک JPEG تنها یا یک تصویر TIFF (شبیه تصاویر و عکس ها) نیست بلکه چند وجهی است: یک صفحه وب گردآوری تعدادی عوامل (متون، تصاویر، نوشته ها، شیوه نامه ها (1)، فایل های صوتی و ویدئویی و غیره) است که ممکن است از جاهایی دیگر (سامانه فایل محلی پایگاه داده وبگاه راه دور دیگر و غیره) بیاید و هنگامی که هر کاربر وبگاه را با یک مرورگر مشاهده می کند، جمع آوری شوند. تحلیل گردآوری یک نمونه 2/9 میلیون وبگاه فرانسوی در سال 2007 نشان داد که حدود 1600 نوع رسانه اینترنتی متفاوت وجود دارد. (2)

• یک وبگاه آغازی ندارد و نمی تواند تا آخر خوانده شود موجودیتی عقلانی است که می تواند توسط کاربری متفاوت از کاربر دیگر دیده شود. در ارتباط با شبکه ای از دیگر وب گاه های پیوند یافته توسط پیوندهای، فرامتنی، شبکه ای از پیوندهایی که حتی قوی تر از شبکه یک انتشارات علمی مبتنی بر استناد ها و اطلاعات کتاب شناختی است وجود دارد.

• وب گاه ها بسیار متعدد هستند و وب هیچ مرزی ندارد. بر اساس آخرین نظر سنجی نتکرافت (3)، بیش از 206 میلیون وبگاه وجود دارد. (4) بیش تر آن ها از هر کشوری در دنیا قابل دسترس هستند و ممکن است بخشی از یک مجموعه ملی بر اساس رهنمودهای خط مشی توسعه مجموعه ملی حقوقی یا سازمانی باشند.

• محتوای وب همیشه درست شبیه یک بخار در حرکت است صفحات وب ممکن است در عرض یک روز چندین مرتبه روزآمد شوند (در صفحات وب روزنامه هایی مانند لوموند (5) و لیبراسیون (6) اعلام کوچکی وجود دارد که نشان می دهند محتوی چند دقیقه قبل روزآمد شده است).

• وب گاه ها و محتوای وب نیز گذرا هستند: یک صفحه وب ممکن است در هر زمانی و به چند دلیل ناپدید شود: توقف اختیاری یا غیر عمدی توسط وب مستر عدم تمدید نام دامنه شکستن دیسک یا مشکلات دسترسی شبکه به سرور میزبان و غیره محتوای وب به یک حادثه خاص پیش بینی شده یا غیر مترقبه پیوند دارد و به ویژه در معرض خطر است به مناسبت یک گزینش مشارکتی و طرح گردآوری وب گاه های سیاسی در طول انتخابات ریاست جمهوری 2007 فرانسه کتابخانه لیون (7) دریافت که 52 درصد از 421 وب گاهی که انتخاب شده بودند یا به طور کامل یا به طور تقریبی پنج ماه بعد از رأی گیری بسته شدند. (8)

چارچوب حقوقی

هر زمان که یک نوع جدید ماده نمایش و ایجاد از جمله فناوری های متنوع جدیدی که در فرانسه ظاهر

ص: 288

style sheets -1

Legal deposit of the French Web: harvesting strategies for a national domain. France Lasfargues, Clément Oury, Bert Wendland, IWAW, 2008: <http://iwaw.net/08/IWAW2008-Lasfargues.pdf> [last accessed on 15-06-2010].

Netcraft -3

Netcraft May 2010 Web Server Survey, <http://news.netcraft.com/archives/category/web-server-survey/> -4

.[[last accessed on 2010-06-15

Le Monde -5

Libé ration -6

Library of Lyon -7

La netcampagne des lé gislatives 2007 en Rhône-Alpes: la course au Net et après,- 8

.[[http://www.pointsdactu.org/article .php3?id_article=863](http://www.pointsdactu.org/article.php3?id_article=863) [last accessed on 2010-06-15

شدند، اختراع میشد کتابخانه ملی فرانسه نخست آزمایش می کرد، سپس سازمانش را با جمع آوری، حفظ و دسترس پذیری به این انتشارات دیجیتال متولد شده وفق می داد زمان آن رسیده است که پس از، کتاب ها حکاکی ها پارتیسیون های موسیقی عکس ها، پوستر ها مدارک صوتی تصویری و چندرسانه ای وب گاه ها نیز بایگانی شوند.

قانون میراث فرانسه (1)، اکنون حق مالکیت IV (ماده 1-311 L تا پایان (1-133 L) از قانون حق مؤلف و نگرش های حقوقی در جامعه اطلاعاتی 2006-961 را که مطابق دستور عمل EC/2001 / 29 پارلمان اروپا و شورای 22 می 2001 در مورد سازگاری جنبه های خاص کپی رایت و حقوق مرتبط با جامعه اطلاعاتی (2) است به ثبت می رساند.

این قانون که به طور رسمی سوم آگوست 2006 منتشر شد:

• دامنه واسپاری حقوقی (3) اینترنت را به این موارد گسترش می دهد: مشمول واسپاری حقوقی هر علامت، نشانه نوشته، تصویر صدا یا پیغام های هر نوع ارتباط با عموم به وسیله کانال های الکترونیکی نیز می شود (ماده 39) قانون برای همه نوع انتشارات الکترونیکی آنلاین از جمله مجموعه ای از علائم نشانه ها، تصاویر، صداها یا هر نوع پیغامی که در اینترنت در دسترس عموم باشد، قابل اجر است نه تنها وب گاه ها، بلکه خبرنامه ها و رسانه های جاری نیز شامل این تعریف می شوند؛

• چگونگی تقسیم مسئولیت های واسپاری وب بین سازمان های تحت قیمیت را تعریف می کند: مؤسسه خبرگزاری ملی که مسئول حفظ میراث صوتی تصویری فرانسه است وب گاه های مرتبط با ارتباطات صوتی تصویری (به طور عمده رادیو و تلویزیون) و کتابخانه ملی فرانسه تمامی وبگاه ها را جمع آوری خواهند کرد؛ یک حکم در دست اقدام است تا فرایندهای گزینش و دسترسی را به اجرا در آورد؛

• راهبردهای جمع آوری را خاص می کند در کتابخانه ملی فرانسه و اسپاری حقوقی اینترنت نیاز به اجازه از ناشران ندارد و به جمع آوری خودکار بخش عمده اولویت می دهد: سازمان های تحت قیمیت ممکن است مواد را مطابق با فرایندهای خاص واسپاری تولید کنندگان از اینترنت جمع آوری کنند. قانون نیز تصریح می کند که هیچ مانعی همچون برقراری ارتباط (4) رمز عبور یا شکل های دیگر محدودیت دسترسی ممکن نیست توسط ناشران برای محدود کردن این فرایند استفاده شود.

دامنه

اگر چه بیشتر وب عمومی توسط هر کس در فرانسه می تواند ملاحظه شود، از نظر فنی و حقوقی غیر ممکن است کل وب را بایگانی کرد. کتابخانه ملی فرانسه دستور دارد تا وب گاه های دامنه ملی فرانسه را جمع آوری کند، یعنی:

ص: 289

The French Heritage Law, or «Code du patrimoine» – 1

Loi n°2006-961 du 1 août 2006 relative au droit d'auteur et aux droits voisins dans la société de – 2

l'information , <http://www.legifrance.gouv.fr/affchTexte.do?cidTexte> = JORF

[TEXT000000266350dateTexte= [last accessed on 2010-06-15

legal depository –3

• به عنوان یک هسته هر وب گاهی با دامنه fr TLD. یا هر TLD مشابه دیگری که به قلمرو مدیریتی فرانسه مربوط می شود (به عنوان مثال re برای جزیره فرانسوی La Reunion)

• هر وب گاهی (به احتمال خارج از دامنه fr) که تولید کننده اش از نظر جغرافیایی تحت قلمرو فرانسه هست (به طور معمول این می تواند در صفحات وب یا با استفاده از سرورهای خاص بررسی شود)؛

• هر وب گاهی (به احتمال خارج از دامنه fr) که بتوان ثابت کرد محتوای تولید شده اش در قلمرو فرانسه به نمایش گذاشته می شود (بررسی این معیار اخیر چالش برانگیز تر است اما فرصتی برای تفسیر و مذاکره با کتابخانه و تولیدکنندگان اینترنت ایجاد می کند).

ابزارها و روشهای جمع آوری

گر چه ما از یک «واسپاری» حقوقی صحبت می کنیم وب گاه ها در حقیقت توسط ناشران به کتابخانه واسپاری نمی شوند. در عوض توسط تکه هایی از نرم افزارهایی به نام رویات های خزنده آرشیو جمع آوری می شوند یک خزنده آرشیو شبیه خزنده های نمایه سازی موتورهای جستجو عمل می کند. برنامه ای است که وب را به روشی خودکار مطابق مجموعه خط مشی هایی مرور می کند. با فهرستی از نشانی های یو.آر.ال (1) شروع و هر صفحه شناسایی شده توسط یو.آر.ال را ذخیره می کند تمامی فرامتن ها را در صفحه می یابد (به عنوان مثال پیوندهایی به دیگر صفحات، تصاویر نوشته ها یا شیوه نامه ها، ویدئوها و غیره) و آن ها را به فهرست یو.آر.ال ها می افزاید تا به طور مسلسل دیده شوند.

پارامترهای فنی بر هویت و رفتار خزنده (دامنه، عمق سرعت فیلترهای تحریم، و غیره) تأثیر می گذارد اما از آن جا که فنون وب بسیار پیچیده هستند و خیلی به سرعت توسعه می یابند، خزنده با بسیاری موانع فنی مواجه می شود که مانع آن از جمع آوری تمامی عوامل یک وب گاه یا حتی یک صفحه وب می شود. بنابراین آرشیوهای وب اغلب ناقص هستند کتابخانه ملی فرانسه از خزنده منبع باز Heritrix که توسط مؤسسات عضو کنسرسیوم حفاظت بین المللی اینترنت توسعه یافته است، (2) استفاده می کند

از آن جا که ممکن نیست جامعیت را هدف قرار داد یا گزینش دستی وب گاه ها را به عهده گرفت، کتابخانه ملی فرانسه در نظر گرفته است دوروش جمع آوری مکمل را برای رویارویی با چالش های واسپاری حقوقی وب ترکیب کند:

• جمع آوری خودکار بخش عمده ای از وب گاه های فرانسوی خزش های وسیع در جهت ارائه که تصویری کلی از هزاران فایل از تعداد بسیار زیادی از وب گاه ها انجام می شود. به عنوان مثال، برای خزش وسیع 2010 که هنوز در زمان نگارش این مقاله هم فعال، است کتابخانه ملی فرانسه بیشترین حد 10000 یو.آر.ال. را برای هر 1/6 میلیون وب گاه جمع آوری می کند خزنده ها محتوی را بدون هیچ تمایزی بین محتوای دانشگاهی سازمانی تجاری یا مبتدل جمع آوری می کنند. این روش به راستی به

ص: 290

URL-1

International Internet Preservation Consortium (IIPC): <http://netpreserve.org> [last accessed on 2010-06-22]

[15] Heritrix crawler: <http://crawler.archive.org> [last accessed on 2010-06-15]

اسلوب واسپاری حقوقی است (به عنوان مثال فرض را بر این قرار نمی دهد آن چه را که مورد علاقه پژوهش گران در 100 سال آینده است، شناسایی کند). با این حال آرشیوهایی که با این روش ایجاد می شوند بسیار سطحی هستند؛ محتوای عمیق یا تحولات و بگاہ را حفظ نمی کنند.

● خزش های متمرکز (1): جمع آوری گزینشی خزش های وسیع را کامل می کند. کتابداران موضوعی وب گاه هایی را برای چنین خزش های متمرکز منطبق با طرح های توسعه مجموعه (همکاری با دیگر کتابخانه ها و پژوهشگران نیز ممکن است) انتخاب می کنند خزش های متمرکز می توانند واقعه محور (انتخابات فرانسه در سال های 2002، 2004، 2007، و 2009) یا موضوعی (خاطرات و بلاگ های شخصی توسعه، پایدار عمل گرای وب و غیره) باشند خزش های متمرکز ساخت آرشیوهای کاملتر و بیشتر را از تعداد محدودی از وب گاه ها ممکن می سازند.

امروزه آرشیوهای وب کتابخانه ملی فرانسه یا آرشیوهای وب حوزه ملی فرانسه شامل 12/5 میلیارد .یو.آر.ال. می شود و 145 ترابایت فضای دیسک اشغال می کند. تاریخ قدیمی ترین صفحات وب به 1996 باز می گردد و این صفحات به لطف آرشیو ملی به دست آمدند که سازمانی غیر انتفاعی است و در نظر دارد یک کتابخانه اینترنتی بسازد و پیشگام آرشیو وب است آخرین تاریخ صفحات وب به چند ساعت قبل باز می گردد.

کشف منبع: خدمات و ابزارهای دسترسی برای کاربران نهایی

تعیین جای خدمت و محدودیت های دسترسی

دسترسی به این آرشیوها مانند دسترسی به اسناد فیزیکی واقع در ساختمان های کتابخانه نیست. وب گاه های گردآوری شده در فهرست کتابخانه ثبت نمی شوند زیرا مجموعه بسیار بزرگ و بسیار ناهمگن است؛ غیر ممکن خواهد بود فهرستی جامع از وب گاه های آرشیو شده ایجاد کرد؛ عنوان ها و محتوای مفصل شان را دانست در، عوض کتابخانه ملی فرانسه فرایندهای نمایه سازی خودکار را ساخته است تا دسترسی سریع به محتوای جمع آوری شده را ممکن سازد. هر فایلی تاریخ می خورد و توصیف می شود تا فقط اطلاعات ضروری (مکان اصلی در وب شکل، اندازه، تعیین جا در آرشیوها، و غیره) را جمع آوری کند. این فرایند نمایه سازی وب گاه های آرشیو شده را قادر می سازد مجدد در محیط انتشاراتی شان نقش داشته باشند و آن ها را با کلیک روی پیوندها درست شبیه وب در حال فعالیت اما در یک محتوای تاریخی و تاریخ خورده مرور کنند.

از آوریل 2008 آرشیوهای وب برای کاربران مجاز در اتاق های مطالعه کتابخانه پژوهشی، در مکان های متفاوت کتابخانه ملی فرانسه (طبقه همکف در کتابخانه فرانسوا میترا و بخش های مجموعه های تخصصی در ریشولیو، لوووا، اوپرا آرسنال و ژان - ویلار در آوینیون) (2) قابل دسترس هستند. گر چه آرشیوها اساساً شامل وب گاه های با دسترسی عمومی و رایگان هستند، این محدودیت وضع

ص: 291

focused crawls -1

Rez-de-jardin level at François-Mitterrand, and special collections depart-ments at Richelieu, Louvois, -2

Opé ra, Arsenal, and Jean-Vilar in Avignon

شده است تا از فراهم آوری های حقوقی که به تمامی مجموعه های واسپاری و میراث حقوقی مربوط می شوند و در نظر دارند کاملاً به کپی رایت و قوانین شخصی احترام بگذارند، تبعیت کند.

برای دسترسی، توافقی کاربران نهایی باید بالاتر از 18 سال باشند و مدرکی جهت نیاز دانشگاهی فعالیت های حرفه ای یا شخصی پژوهشی شان برای دسترسی به این آرشیوها ارائه کنند (اکنون کتابخانه ملی فرانسه آخرین و تنها منبع این نوع سند است زیرا تنها کتابخانه ای است که این خدمت را در فرانسه ارائه می کند). کارت های خوانندگان توسط سرویس راهنمای خوانندگان در کتابخانه های فرانسوا-میتران یا ریشولیو پس از مصاحبه شخصی با یک کتابدار صادر می شود بر اساس نیازهای کاربران این مصاحبه تعیین می کند کاربران اجازه ورود به یک یا چند بخش را دارند و دوره اعتبار کارت (3 روز 15 روز سالانه) چقدر است

آرشیوهای وب همراه با تمامی خدمات الکترونیکی کتابخانه (وب گاه های اطلاعاتی، تسهیلات رزرو، فهرست ها، دسترسی به اینترنت) و منابع الکترونیکی (کتاب ها و تصاویر دیجیتالی، نشریات الکترونیکی پایگاه داده های آنلاین، لوح های فشرده، بوک مارک ها) روی 350 کامپیوتر واقع در اتاق های مطالعه کتابخانه پژوهشی در دسترس هستند این کامپیوترها در دسترس عموم هستند اما کاربران برای یک جا در کتابخانه و استفاده از کامپیوترها نیاز دارند از قبل جا ذخیره کنند

ابزارهای جستجو و دیدن

آرشیوهای وب از طریق برنامه کاربردی اختصاصی ای که مجموعه «آرشیوهای اینترنت» (1) نامیده می شود، قابل دسترسی می شوند و از «آرشیوها» استفاده می شود تا تأکید کند که این مجموعه ها جامع نیستند. برنامه کاربردی توسط یک نوار نرنجی با حروف WWW و یک نشانگر ماوس روی آن ارائه می شود تا نشان دهد گر چه صفحات وب در یک جعبه قرار می گیرند هنوز قابل کلیک هستند کتابخانه ملی فرانسه توسعه یک همانندی دیداری و یک علامت خاص (2) را انتخاب کرده است تا آرشیوهای وب را بیش تر قابل دیدن کند.

عکس

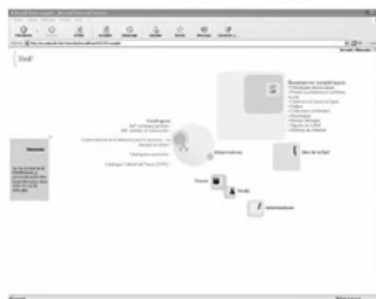
شده است تا از فراهم‌آوری‌های حقوقی که به تمامی مجموعه‌های واسپاری و میراث حقوقی مربوط می‌شوند و در نظر دارند کاملاً به کپی‌رایت و قوانین شخصی احترام بگذارند، تبعیت کند.

برای دسترسی توافقی، کاربران نهایی باید بالاتر از ۱۸ سال باشند و مدرکی جهت نیاز دانشگاهی، فعالیت‌های حرفه‌ای، یا شخصی پژوهشی‌شان برای دسترسی به این آرشیوها ارائه کنند (اکنون کتابخانه ملی فرانسه آخرین و تنها منبع این نوع سند است، زیرا تنها کتابخانه‌ای است که این خدمت را در فرانسه ارائه می‌کند). کارت‌های خوانندگان توسط سرویس راهنمای خوانندگان در کتابخانه‌های فرانسوا-میتران یا ریشولیو پس از مصاحبه شخصی با یک کتابدار صادر می‌شود. براساس نیازهای کاربران، این مصاحبه تعیین می‌کند کاربران اجازه ورود به یک یا چند بخش را دارند و دوره اعتبار کارت (۳ روز، ۱۵ روز، سالانه) چقدر است.

آرشیوهای وب همراه با تمامی خدمات الکترونیکی کتابخانه (وبگاه‌های اطلاعاتی، تسهیلات رزرو، فهرست‌ها، دسترسی به اینترنت) و منابع الکترونیکی (کتاب‌ها و تصاویر دیجیتالی، نشریات الکترونیکی، پایگاه‌داده‌های آنلاین، لوح‌های فشرده، بوک‌مارک‌ها) روی ۳۵۰ کامپیوتر واقع در اتاق‌های مطالعه کتابخانه پژوهشی در دسترس هستند. این کامپیوترها در دسترس عموم هستند، اما کاربران برای یک جا در کتابخانه و استفاده از کامپیوترها نیاز دارند از قبل جا ذخیره کنند.

ابزارهای جستجو و دیدن

آرشیوهای وب از طریق برنامه کاربردی اختصاصی‌ای که مجموعه «آرشیوهای اینترنت»^۱ نامیده می‌شود، قابل دسترس می‌شوند. واژه «آرشیوها» استفاده می‌شود تا تأکید کند که این مجموعه‌ها جامع نیستند. برنامه کاربردی توسط یک نواز نارنجی با حروف www و یک نشانگر ماوس روی آن ارائه می‌شود تا نشان دهد گرچه صفحات وب در یک جعبه قرار می‌گیرند، هنوز قابل کلیک هستند. کتابخانه ملی فرانسه توسعه یک همانندی دیداری و یک علامت خاص^۲ را انتخاب کرده است تا آرشیوهای وب را بیشتر قابل دیدن کند.



تصویر ۱. صفحه اصلی خدمات و منابع الکترونیکی کتابخانه پژوهشی کتابخانه ملی فرانسه

1. Archives de l'Internet
2. logo

تصویر ۱. صفحه اصلی خدمات و منابع الکترونیکی کتابخانه پژوهشی کتابخانه ملی فرانسه

برای مرور آرشیوها کتابخانه ملی فرانسه سه ابزار متفاوت ارائه می کند:

- جستجو با یو آر. ال.
 - جستجو با کلیدواژه
 - مجموعه های مشخصی که در نظر دارند کشف مجموعه های اختصاصی را تسهیل کنند.
- هر سه ابزار در برنامه کاربردی واحدی به نام آرشیو اینترنت منسجم شده است (تصویر 2)

عکس

برای مرور آرشیوها، کتابخانه ملی فرانسه سه ابزار متفاوت ارائه می‌کند:

- جستجو با یو.آر.آل.
 - جستجو با کلیدواژه.
 - مجموعه‌های مشخصی که در نظر دارند کشف مجموعه‌های اختصاصی را تسهیل کنند.
- هر سه ابزار در برنامه کاربردی واحدی به نام «آرشیو اینترنت» منسجم شده است (تصویر ۲).



تصویر ۲. صفحه خانگی عملکرد آرشیو اینترنت

- جستجوی یو.آر.آل.^۱ کاربران را قادر می‌سازد آرشیو یک وبگاه، یک صفحه وب یا حتی فایل را با وارد کردن جایگاه اصلی اینترنتی‌اش جستجو کنند. به عنوان مثال، جستجوی <http://www.lemonde.fr> ۶۱۷ نتیجه می‌دهد (از ۱۵ ژوئن ۲۰۱۰) که در یک نگاه تقویمی از ۱۹۹۶ تا ۲۰۱۰ نمایش داده می‌شوند (نتایج می‌توانند به تاریخ‌های خاصی محدود شوند). هر تاریخ قابل کلیک است و دسترسی به صفحه خانگی روزنامه لوموند می‌دهد که در آن تاریخ قابل دسترس بوده است. این نوع جستجو زمان جستجوی توسعه یک وبگاه و مقایسه ویرایش‌های آن بسیار مفید است، اما درست شبیه جستجوی یک شماره بازیابی در یک فهرست، کاربران باید یو.آر.آل. وبگاه یا صفحه وبی که می‌خواهند جستجو کنند، بدانند (به عنوان مثال، از یک کتابشناسی) یا باید قادر باشند آن را پیدا کنند (با استفاده از راهنماها یا پیوندهای بیرونی^۲ که به وبگاه‌های همکار ارجاع می‌دهند). کاربران نیز ممکن است وب در حال فعالیت و آرشیوهای وب را هم‌زمان در حال استفاده از دو مرورگر با ظاهری متفاوت، به نوبت مرور کنند. در حال حاضر، جستجوی یو.آر.آل. تنها روش جستجوی تمام آرشیوهاست.
- جستجوی کلیدواژه‌ای شبیه یک موتور جستجوی سنتی کار می‌کند: کاربران را قادر می‌سازد اسناد متنی را که شامل یک یا چند واژه است، جستجو نمایند. به لطف گزینه‌های جستجوی پیشرفته نیز کاربران ممکن است عین یک عبارت یا جمله را جستجو کنند یا جستجو را به یک وبگاه خاص

۱. URL کونه‌نوشت Uniform Resource Locator است و اشاره به جایگاه اینترنتی صفحه وب دارد.

2. out-links

تصویر 2. صفحه خانگی عملکرد آرشیو اینترنت

جستجوی یو.آر.آل. (1) کاربران را قادر می‌سازد آرشیو یک وبگاه یک صفحه وب یا حتی فایل را با وارد کردن جایگاه اصلی اینترنتی‌اش جستجو کنند. به عنوان مثال جستجوی <http://www.lemonde.fr> 617 نتیجه می‌دهد (از 15 ژوئن 2010) که در یک نگاه تقویمی از 1996 تا 2010 نمایش داده می‌شوند (نتایج می‌توانند به تاریخ‌های خاصی محدود شوند). هر تاریخ قابل کلیک است و دسترسی به صفحه خانگی روزنامه لوموند می‌دهد که در آن تاریخ قابل دسترس بوده است. این نوع جستجو زمان جستجوی توسعه یک وبگاه و مقایسه ویرایش‌های آن بسیار مفید است اما درست شبیه جستجوی یک شماره بازیابی در یک فهرست، کاربران باید یو.آر.آل. وب‌گاه یا

صفحه وبی که می خواهند جستجو کنند، بدانند (به عنوان مثال از یک کتابشناسی) یا باید قادر باشند آن را پیدا کنند (با استفاده از راهنماها یا پیوندهای بیرونی (2) که به وب گاه های همکار ارجاع می دهند). کاربران نیز ممکن است وب در حال فعالیت و آرشیوهای وب را همزمان در حال استفاده از دو مرورگر با ظاهری متفاوت به نوبت مرور کنند. در حال حاضر، جستجوی یو. آر. ال. تنها روش جستجوی تمام آرشیوهاست.

• جستجوی کلید واژه ای شبیه یک موتور جستجوی سنتی کار می کند: کاربران را قادر می سازد اسناد متنی را که شامل یک یا چند واژه است جستجو نمایند. به لطف گزینه های جستجوی پیشرفته نیز کاربران ممکن است عین یک عبارت یا جمله را جستجو کنند یا جستجو را به یک وب گاه خاص

ص: 293

1- URL کوتاه نوشت Uniform Resource Locator است و اشاره به جایگاه اینترنتی صفحه وب دارد.

2- out-links

محدود کنند. این جستجو هنوز آزمایشی است و فقط پنج درصد آرشیو های وب کتابخانه ملی فرانسه را شامل می شود نمایه سازی تمام متن میلیاردها فایل بینظم با به حساب آوردن نسخه های تکراری و انسجام موقت که آن ها را به یکدیگر پیوند می دهد، یک چالش فنی است که چندین طرح پژوهشی بین المللی در تلاش اند تا آن را برطرف کنند اما هنوز موفق نشده اند.

• علاوه بر گزینه های جستجو مهمترین کارآمدی های این عملکرد عبارت اند از: توانایی نمایش صفحات وب با جمع آوری تعدادی عوامل که ممکن است در تاریخ های متفاوتی یا دست کم در زمان های مختلف آرشیو شده باشند و در نتیجه دوباره یک وبگاه مصنوعی ایجاد کنند.

• توانایی کلیک روی پیوندها با در نظر گرفتن انسجام موقت (به عنوان مثال اگر ما به وبگاه الیزه ریاست جمهوری فرانسه (1) در 7 می 2007 نگاه کنیم و روی یک پیوند کلیک کنیم تا وبگاه دولت را ببینیم انتظار داریم ویرایش ماه می 2007 را مشاهده کنیم).

این کارآمدی ها و جستجوی یو.آر.آل توسط ماشین وی یک منبع باز (2) پشتیبانی می شوند و توسط آرشیو اینترنت با مشارکت ICP توسعه می یابند جستجوی کلیدواژه ای بر پایه نرم افزار ناچ واکس منبع باز (3) (ناچ با الحاقات آرشیو وب) بنا شده است. همچنین این نرم افزار که توسط آرشیو اینترنت و گردهمایی کتابخانه های ملی نوردیک (4) توسعه یافته ابزاری است برای نمایه سازی و جستجوی آرشیوهای وب با استفاده از موتور جستجوی ناچ و الحاقات که آرشیوهای وب را جستجو می کند. کتابخانه ملی فرانسه جهت ساخت یک برنامه کاربردی بر اساس این ابزارها آن ها را در ابزارها و فرایندهای توسعه اش یکپارچه و سفارشی کرده و به کارآمدی های کوچک تر گسترش داده است (به عنوان مثال، اجرای طرحی برای سؤالات مداوم یو.آر.آل.).

مجموعه های ویژه

کتابداران موضوعی کتابخانه ملی فرانسه جهت جبران ابهام مجموعه ها و فقدان ابزارهای دسترسی (که هنوز توسعه می یابند) و نیز آگاهی در مورد مجموعه های موضوعی غنی، با همکاری پژوهشگران، مسیرهای هدایت شده (به معنی واقعی کلمه: تورهای هدایت شده یا مجموعه های ویژه در فرانسه) را ساخته اند که مجموعه ای از صفحات مصور و منظم ویرایشی است که شامل پیوندهای مستقیم به آرشیوهای وب می شود مقدمه ای ارائه می کند جستجو را شبیه سازی می نماید و ایده هایی در مورد دیگر جستجوهای ممکن می دهد از سال 2008 سه مجموعه ویژه منتشر شده است:

• کلیک کن و رأی بده: وب انتخاباتی: گزینشی از وب گاه ها از انواع تهیه کننده (سازمان ها، کاندیداها، حامیان، تماشاچیان، افراد و غیره) از مبارزات انتخاباتی، 2002، 2004 و 2007 تصویر (3)؛

ص: 294

the French Presidency Ellysé e website – 1

The open source Wayback Machine, <http://archive-access.sourceforge.net/projects/wayback/> [last – 2
[accessed on 2010-06-15]

NutchWAX (Nutch Web Archive eXtensions), <http://archive-access.sourceforge.net/projects/nutch/> – 3
.[[last accessed on 2010-06-15]

● در مورد خود نوشتن در وب خاطرات شخصی و ادبی: در نظر دارد نتایج انتقال از کاغذ را به وب و روشی که بلاگ ها نوشته های شخصی ادبی و انتقادی را تغییر داده اند، بررسی کند؛

● عمل گرایی وب (1): نشان می دهد چطور پس از سال ها وب توسط عمل گراها به عنوان ابزار انتشار و، ارتباط ابزاری که قبول مسئولیت را تشویق می کند و مکانی برای، بحث مشارکت سازماندهی و عمل استفاده شده است.

عکس

معرفی آرشیوهای وب ... ۲۹۵

- در مورد خود نوشتن در وب: خاطرات شخصی و ادبی: در نظر دارد نتایج انتقال از کاغذ را به وب و روشی که بلاگ ها نوشته های شخصی، ادبی، و انتقادی را تغییر داده اند، بررسی کند؛
- عمل گرایی وب: نشان می دهد چطور پس از سال ها، وب توسط عمل گراها به عنوان ابزار انتشار و ارتباط ابزاری که قبول مسئولیت را تشویق می کند و مکانی برای بحث، مشارکت، سازماندهی، و عمل استفاده شده است.



تصویر ۳. مجموعه ویژه «کلیک کن و رأی بده: وب انتخاباتی»

زمانی که این مقاله نوشته می شود، سه مجموعه ویژه تر در دست تهیه است: توسعه پایدار، فیلم های غیر حرفه ای، و سفرنامه نویسی.

به عنوان یک خدمت مکمل، کاربران اجازه دارند صفحات وب را چاپ و کپی کنند و نمونه های متن را در یک فایل نوت پد^۱ درج کنند، محتوای نوت پد را چاپ کنند یا آن را از طریق ایمیل بفرستند. به علت محدودیت های حقوقی، در حال حاضر تسهیلاتی وجود ندارد که کاربران را قادر سازد از صفحه عکس بگیرند یا عناصر وب را از آرشیوها بردارند.

طرح های توسعه آتی، نمایه سازی تمام متن و جستجو به علاوه دسترسی ناپیوسته و بگانه از کتابخانه های منطقه ای را همراه با تسهیم مأموریت و اسپاری حقوقی با کتابخانه ملی فرانسه در اطراف کشور دربرمی گیرد.

1. Web activism
2. note pad

تصویر 3. مجموعه ویژه «کلیک کن و رأی بده: وب انتخاباتی»

زمانی که این مقاله نوشته می شود سه مجموعه ویژه تر در دست تهیه است توسعه پایدار فیلم های غیر حرفه ای، و سفرنامه نویسی.

به عنوان یک خدمت مگمل، کاربران اجازه دارند صفحات وب را چاپ و کپی کنند و نمونه های متن را در یک فایل نوت پد (2) درج کنند محتوای نوت پد را چاپ کنند یا آن را از طریق ایمیل بفرستند. به علت محدودیت های حقوقی در حال حاضر تسهیلاتی وجود ندارد که کاربران را قادر سازد از صفحه عکس بگیرند یا عناصر وب را از آرشیوها بردارند.

طرح های توسعه آتی نمایه سازی تمام متن و جستجو به علاوه دسترسی ناپیوسته وب گاه از کتابخانه های منطقه ای را همراه با تسهیم مأموریت و اسپاری حقوقی با کتابخانه ملی فرانسه در اطراف کشور در بر می گیرد.

ص: 295

Web activism -1

note pad -2

ارزش آرشیوهای وب و علاقه به آن ها فقط با گذشت زمان ثابت می شود پس از اینکه منابع وب از وب ناپدید شدند؛ این بخشی از فلسفه میراث کتابخانه ملی فرانسه و سنتی نیست که انتظار یک بازگشت سرمایه کوتاه مدت را داشته باشد؛ رسالتش فهم اهمیت زمان است.

با این حال نخستین ارزیابی کاربران و کاربرد آرشیوهای وب مرحله ای مهم در اجرای واسپاری حقوقی اینترنت است. اثبات سودمندی عمومی و علمی مجموعه های جمع آوری شده و نیز توسعه و تحلیل نخستین آمارهای کاربرد ما را قادر خواهد ساخت توان بالقوه این مجموعه ها و نیز محدودیت های شان را اندازه گیری کنیم این تحلیل ها در رویارویی با انتظارات پژوهشگران برای توسعه مجموعه و ابزار هر دو مفید خواهند بود

تحلیل کمی

ابزاری تحلیلی به نام AWStats ایجاد شده است تا ترافیک همزمان عملکرد را تحلیل کند. شبیه ماشین وی بک و ناچ واکس درون سازمانی سفارشی ساخته شده است تا بین استفاده عمومی توسط خوانندگان، استفاده مرجع توسط کتابداران در هنگام ورود یا در میز مرجع و استفاده حرفه ای توسط کتابداران موضوعی کتابخانه ملی فرانسه تمایز ایجاد نماید

عکس

کاربرد: اطلاعات و ارقام

ارزش آرشیوهای وب و علاقه به آنها فقط با گذشت زمان ثابت می‌شود، پس از اینکه منابع وب از وب ناپدید شدند؛ این بخشی از فلسفه میراث کتابخانه ملی فرانسه و سستی نیست که انتظار یک بازگشت سرمایه کوتاه‌مدت را داشته باشد؛ رسالتش فهم اهمیت زمان است.

با این حال، نخستین ارزیابی کاربران و کاربرد آرشیوهای وب مرحله‌ای مهم در اجرای واسپاری حقوقی اینترنت است. اثبات سودمندی عمومی و علمی مجموعه‌های جمع‌آوری شده و نیز توسعه و تحلیل نخستین آمارهای کاربرد، ما را قادر خواهد ساخت توان بالقوه این مجموعه‌ها و نیز محدودیت‌هایشان را اندازه‌گیری کنیم. این تحلیل‌ها در رویارویی با انتظارات پژوهشگران برای توسعه مجموعه و ابزار هر دو مفید خواهند بود.

تحلیل کمی

ابزاری تحلیلی به نام AWStats ایجاد شده است تا ترافیک هم‌زمان عملکرد را تحلیل کند. شبیه ماشین وی‌بک و ناچ‌واکس، درون‌سازمانی سفارشی ساخته شده است تا بین استفاده عمومی توسط خوانندگان، استفاده مرجع توسط کتابداران در هنگام ورود یا در میز مرجع، و استفاده حرفه‌ای توسط کتابداران موضوعی کتابخانه ملی فرانسه تمایز ایجاد نماید.

جدول ۱. شاخص‌های کاربرد اصلی در سال ۲۰۰۸ و ۲۰۰۹

۲۰۰۹	۲۰۰۸	
۱۰۶	۳۵	میانگین تعداد جلسه‌ها در هر ماه
۱۲۷۵	۳۱۶	تعداد کل جلسه‌ها
۹۰۰۶۳	۳۵۸۹۱	تعداد کل صفحات مشاهده شده

جدول ۲. کمیت‌های جامع در سال ۲۰۰۹

صفحات مشاهده شده	جلسه‌های طولانی (مشاهده‌های بیشتر از ۱ ساعت)	جلسه‌ها (مشاهده‌های بیشتر از ۵ دقیقه)	مشاهده‌کنندگان	مشاهده‌ها	
۱۴۰۴۵	۷۱	۳۳۸	۴۸۱	۹۰۴	عمومی
۶۲۴۲	۲۳	۹۴	۱۱۸	۲۱۶	مرجع
۶۹۷۷۶	۲۳۹	۸۴۳	۵۵۲	۱۹۴۵	حرفه‌ای
۹۰۰۶۳	۳۳۳	۱۲۷۵	۱۲۳۵	۳۰۶۵	کل

جدول ۱. شاخص‌های کاربرد اصلی در سال ۲۰۰۸ و ۲۰۰۹

جدول ۲. کمیت‌های جامع در سال ۲۰۰۹

دو نتیجه ای که توجه به کتابخانه ملی فرانسه را جلب می کند عبارت اند از:

• تعداد کل جلسه ها بین سال 2008 و 2009 سه بار افزایش یافته است (بعدها خواهیم دید که آموزش اصلی و ابتکار عمل های اطلاعاتی برای خوانندگان و کتابداران مرجع توسعه یافته است)؛

• تعداد فزاینده ای از کاربران نهایی وجود دارند که آرشیوها را برای پژوهش عمیق استفاده می کنند. میانگین جلسه در فوریه، 2009، 13 دقیقه طول کشید و در دسامبر 2009 تا 30 دقیقه افزایش یافت بین فوریه و سپتامبر 2009 بیش از دو جلسه وجود نداشت که بیش از یک ساعت یا بیشتر طول بکشد. سپس، تعداد هر جلسه در اکتبر 2009 به 7 جلسه در نوامبر به 32 جلسه، در دسامبر به 22، جلسه و در ژانویه 2010 به 23 جلسه افزایش یافت این جلسه ها نشان می دهد که اجرای پژوهش وسیع روی آرشیوهای وب آغاز می شود با نگاهی به فهرست بیش ترین صفحات مشاهده شده به نظر می رسد این طرح ها بیش از همه توسط پژوهشگران علوم اجتماعی و علوم سیاسی اجرا می شود.

تحلیل کیفی

در تکمیل تحلیل آماری و فهرستی جامع از یو.آر.ال.ها و کلید واژه های جست و جو شده، کتابخانه ملی فرانسه یک لاگ داخلی و مشترک مورد استفاده کتابدارانی ایجاد کرد که کاربران را به مناسبت های مختلف از جمله مصاحبه پذیرش میز، مرجع درخواست ملاقات، شخصی، نمایش به تهیه کنندگان وب گاه ها و غیره ملاقات می کنند از آوریل 2008 کتابداران 34 کاربر نهایی را با وارد کردن اطلاعاتی همچون تاریخ/ساعت گروه و نام کتابدار شکل در خواست (در محل، تلفنی، ایمیل، ...)، نام نوع، کارت موضوع پژوهش، یادداشت ها /سؤال ها نظرات/ اظهارات، ثبت نام کردند (تصویر 4).

عکس

دو نتیجه‌ای که توجه به کتابخانه ملی فرانسه را جلب می‌کند عبارت‌اند از:

- تعداد کل جلسه‌ها بین سال ۲۰۰۸ و ۲۰۰۹ سه بار افزایش یافته است (بعدها خواهیم دید که آموزش اصلی و ابتکار عمل‌های اطلاعاتی برای خوانندگان و کتابداران مرجع توسعه یافته است)؛
- تعداد فزاینده‌ای از کاربران نهایی وجود دارند که آرشیوها را برای پژوهش عمیق استفاده می‌کنند. میانگین جلسه در فوریه ۲۰۰۹، ۱۳ دقیقه طول کشید و در دسامبر ۲۰۰۹ تا ۳۰ دقیقه افزایش یافت. بین فوریه و سپتامبر ۲۰۰۹، بیش از دو جلسه وجود نداشت که بیش از یک ساعت یا بیشتر طول بکشد. سپس، تعداد هر جلسه در اکتبر ۲۰۰۹ به ۷ جلسه، در نوامبر به ۳۲ جلسه، در دسامبر به ۲۲ جلسه، و در ژانویه ۲۰۱۰ به ۲۳ جلسه افزایش یافت. این جلسه‌ها نشان می‌دهد که اجرای پژوهش وسیع روی آرشیوهای وب آغاز می‌شود. با نگاهی به فهرست بیشترین صفحات مشاهده شده، به‌منظر می‌رسد این طرح‌ها بیش از همه توسط پژوهشگران علوم اجتماعی و علوم سیاسی اجرا می‌شود.

تحلیل کیفی

در تکمیل تحلیل آماری و فهرستی جامع از یو.آر.ال‌ها و کلیدواژه‌های جست‌وجوشده، کتابخانه ملی فرانسه یک لاگ داخلی و مشترک مورد استفاده کتابدارانی ایجاد کرد که کاربران را به مناسبتهای مختلف از جمله مصاحبه پذیرش، میز مرجع، درخواست ملاقات شخصی، نمایش به تهیه‌کنندگان وبگاه‌ها، و غیره ملاقات می‌کنند. از آوریل ۲۰۰۸، کتابداران ۳۴ کاربر نهایی را با وارد کردن اطلاعاتی همچون تاریخ/ساعت، گروه و نام کتابدار، شکل درخواست (در محل، تلفنی، ایمیل، ...)، نام، نوع کارت، موضوع پژوهش، یادداشت‌ها/سؤال‌ها/نظرات/اظهارات، ثبت‌نام کردند (تصویر ۴).

سوالات	پاسخ‌ها
تاریخ و ساعت	می / ژانویه / اوت ۲۰۰۹
بخش یا خدمت و نام کارگزار	کریستین ژن
پشتیبانی از درخواست	در محل
نام کوچک و نام خانوادگی خواننده	
عنوان دسترسی	نمایش‌ها
موضوع پژوهش / رشته	<p>۲۹ می (۱۷ ساعت و نیم): الیزابت لگرو، بلاگر و عضو گروه یار APA http://2009sediments.wordpress.com/</p> <p>۲۰ ژانویه (۱۶ ساعت): مارتین سونه، پژوهشگر CNRS، نویسنده و بلاگر www.martinesonnet.fr/blogwp</p> <p>۱۱ اوت (۱۷ ساعت): سیلور مرسیه که یادداشتی را در بلاگش نوشته است: http://www.bibliosession.net/2009/09/17/archives-de-liternet-demadez-votre-ticket-pour-la-posterite/</p> <p>۲۱ اوت (۱۷ ساعت): اوریان دزینی که در Paris XIII تدریس می‌کند و نویسنده رساله‌ای در مورد روزنامه‌های خصوصی آنلاین وی به‌ویژه علاقه‌مند به مسیر هدایت‌شده است و از یک همکار بانک سالن در مورد این موضوع سؤال کرده است که آن را به‌سوی من ارجاع داده است.</p>
مشاهدات	

سؤالات	پاسخها
تاریخ و ساعت	۲۴ نوامبر ۲۰۰۹
بخش یا خدمت و نام کارگزار	SOL کریستوف تربویی
پشتیبانی از درخواست	در یک محل رسمی
نام کوچک و نام خانوادگی خواننده	فابین گرفه
عنوان دسترسی	کارت سالانه (۲۰۰۹/۱۱/۲۴ تحویل شد)
موضوع پژوهش / رشته	مبارزه انتخاباتی اروپایی سال ۲۰۰۹ در اینترنت
مشاهدات	

سؤالات	پاسخها
تاریخ و ساعت	۲۰۰۹/۸/۱۹
بخش یا خدمت و نام کارگزار	SOL
پشتیبانی از درخواست	در محل
نام کوچک و نام خانوادگی خواننده	کارول دافینی
عنوان دسترسی	کارت سالانه
موضوع پژوهش / رشته	آرشیوهای وب (کارآموز محافظه کار)
مشاهدات	

تصویر ۴. لاگ کاربران نهایی آرشیوهای وب

این ابزار ساده، کتابخانه ملی فرانسه را قادر ساخت مجموعه‌هایی را که مورد درخواست هستند، شناسایی کند؛ این مجموعه‌ها در حال حاضر در حوزه علوم اجتماعی و علوم سیاسی هستند. روزی که کتابخانه ملی فرانسه خدمت را عرضه کرد، یک دانشجوی کارشناسی ارشد که در مورد «اینترنت و انتخابات ریاست جمهوری ۲۰۰۷» کار می‌کرد، درخواست دسترسی کرد. یک دانشجوی دکترای زبان‌شناسی که در مورد تحلیل سخنرانی‌های کاندیداهای زن کار می‌کرد از روم، ایتالیا آمده بود تا در مورد بلاگ رهبر کمونیست، ماری - ژرژ بوفه^۱ که بعد از انتخابات ریاست جمهوری بسته شده بود، کار کند. فعالیت وب نیکلا سارکوزی، حزب سوسیالیست، انتخابات اروپایی ۲۰۰۹، استفاده از ویدئوها در فعالیت انتخاباتی، کارتونها و کاریکاتورهای سیاسی، نمونه‌هایی از پژوهش در این حوزه‌ها هستند.

پژوهش‌های دیگر شامل موضوعاتی بوده است همچون جستجوی اطلاعات در اتحادیه اروپا، آرشیوهای وزارت بوم‌شناسی و موجودیت‌های نامتمرکزش، وب‌سایت‌های نویسندگان غیر حرفه‌ای، وب‌سایت‌های مدیریت استرس، وب‌سایت‌های شخصی و سؤالاتی درباره اینکه چرا و چگونه

1. Marie-George Buffet

تصویر ۴. لاگ کاربران نهایی آرشیوهای وب

این ابزار ساده کتابخانه ملی فرانسه را قادر ساخت مجموعه‌هایی را که مورد درخواست هستند شناسایی کند؛ این مجموعه‌ها در حال حاضر در حوزه علوم اجتماعی و علوم سیاسی هستند. روزی که کتابخانه ملی فرانسه خدمت را عرضه کرد یک دانشجوی کارشناسی ارشد که در مورد «اینترنت و انتخابات ریاست جمهوری ۲۰۰۷» کار می‌کرد درخواست دسترسی کرد. یک دانشجوی دکترای زبان‌شناسی که

در مورد تحلیل سخنرانی های کاندیداهای زن کار می کرد از، روم ایتالیا آمده بود تا در مورد بلاگ رهبر کمونیست ماری - ژرژ بوفه (1) که بعد از انتخابات ریاست جمهوری بسته شده بود، کار کند. فعالیت وب نیکلا سارکوزی، حزب سوسیالیست انتخابات اروپایی 2009 استفاده از ویدئوها در فعالیت انتخاباتی، کارتون ها و کاریکاتورهای، سیاسی نمونه هایی از پژوهش در این حوزه ها هستند.

پژوهش های دیگر شامل موضوعاتی بوده است همچون جستجوی اطلاعات در اتحادیه اروپا آرشیوهای وزارت بوم شناسی و موجودیت های نامتمرکزش وب سایت های نویسندگان غیر حرفه ای وب سایت های مدیریت استرس وب سایت های شخصی و سؤالاتی درباره این که چرا و چگونه

ص: 298

وب سایت ها جمع آوری می شوند و قوانین رقابتی تا در یک مورد حقوقی کمک کند.

نظر سنجی ها

نخستین نظر سنجی از اکتبر 2006 تا ژوئن 2007 در محتوای یک دوره درسی در مقطع کارشناسی ارشد به نام «اینترنت در طول فعالیت» انجام شد و توسط کتابخانه ملی فرانسه و یک پژوهشگر دانشگاه در رسانه اجتماعی به طور هماهنگ سازماندهی شد. 17 جلسه مشاهده و 5 مصاحبه انجام شد تا نیازهای کاربران به ابزارها و کارآمدی ها و نگرش شان به یک نوع رسانه جدید شناسایی شوند (برای آزمون نهایی آن ها، دانشجویان مجبور شدند مقاله ای بنویسند در موضوعی که مربوط به منبع وب در حال فعالیت و محتوای آرشیو شده اش بود).

دومین نظر سنجی برای نوامبر 2010 برنامه ریزی می شود این نظر سنجی کاربران فعلی و بالقوه را در نظر دارد و قصد دارد نیازهای شان را به محتوای مجموعه تعریف و ابزار توسعه کاربرد را شناسایی کند (جدول 3).

عکس

وبسایت‌ها جمع‌آوری می‌شوند و قوانین رقابتی تا در یک مورد حقوقی کمک کند.

نظرسنجی‌ها

نخستین نظرسنجی از اکتبر ۲۰۰۶ تا ژوئن ۲۰۰۷ در محتوای یک دوره درسی در مقطع کارشناسی ارشد به نام «اینترنت در طول فعالیت» انجام شد و توسط کتابخانه ملی فرانسه و یک پژوهشگر دانشگاه در رسانه اجتماعی به‌طور هماهنگ سازماندهی شد. ۱۷ جلسه مشاهده و ۵ مصاحبه انجام شد تا نیازهای کاربران به ابزارها و کارآمدی‌ها و نگرش‌شان به یک نوع رسانه جدید شناسایی شوند (برای آزمون نهایی آنها، دانشجویان مجبور شدند مقاله‌ای بنویسند در موضوعی که مربوط به منبع وب در حال فعالیت و محتوای آرشیو شده‌اش بود).

دومین نظرسنجی برای نوامبر ۲۰۱۰ برنامه‌ریزی می‌شود. این نظرسنجی کاربران فعلی و بالقوه را در نظر دارد و قصد دارد نیازهایشان را به محتوای مجموعه تعریف و ابزار توسعه کاربرد را شناسایی کند (جدول ۳).

جدول ۳. چارچوبی برای نظرسنجی بعدی کاربرد آرشیو وب در نوامبر ۲۰۱۰

کاربران بالقوه آرشیو وب	کاربران فعلی آرشیو وب
<p>۵-۸ مصاحبه برای هر گروه هدف:</p> <ul style="list-style-type: none"> • دانشمندان در دیگر حوزه‌های پژوهشی (هنرهای دیجیتال، علوم انسانی، پژوهشگران رسانه‌ای علاقه‌مند به خود وب، و غیره). • یک گروه وسیع‌تر علاقه‌مند به حافظه وب، وکلا و متخصصان اطلاعات. 	<p>۵-۸ مصاحبه</p>
<p>سؤالاتی درباره:</p> <ul style="list-style-type: none"> • آگاهی از آرشیوهای وب • علاقه شناسایی شده به خود حوزه پژوهشی • انواع اطلاعات جستجو شده اگر مورد نیاز هستند • ضروریات کاربرد واقعی 	<p>سؤالاتی درباره:</p> <ul style="list-style-type: none"> • خط‌مشی توسعه مجموعه: گزینش، تناوب و کیفیت جمع‌آوری • کاربران هدف و انواع استفاده برای توسعه • علاقه و موضوعاتی برای مجموعه‌های ویژه • اطلاعات و متادیتاهای ارائه شده در نمایش نتایج و نمایش یک صفحه

مشارکت: راهبردهایی برای تأسیس یک انجمن کتابداران و رسیدن به کاربران نهایی

معرفی آرشیوهای وب در کتابخانه به معنی ساخت یک مفهوم کلی از مالکیت این نوع جدید مجموعه توسط هر دو کارکنان و کاربران نهایی است.

مجموعه‌ها به شیوه‌های مختلفی باعث سرخوردگی می‌شوند (ذخیره‌های وب ناقص هستند و برخی وبگاه‌ها در آرشیوها هم‌زمان ظاهر نمی‌شوند). مانند سرخوردگی عملکرد جستجو: برای کارکنان، به‌ویژه، فقدان یک فهرست و این حقیقت که کتابداران هر کار نظام‌مند توصیفی را در مورد وبگاه‌ها انجام

جدول 3. چارچوبی برای نظرسنجی بعدی کاربرد آرشیو وب در نوامبر 2010

مشارکت: راهبردهایی برای تأسیس یک انجمن کتابداران و رسیدن به کاربران نهایی

معرفی آرشیوهای وب در کتابخانه به معنی ساخت یک مفهوم کلی از مالکیت این نوع جدید مجموعه توسط هر دو کارکنان و کاربران نهایی است.

مجموعه‌ها به شیوه‌های مختلفی باعث سرخوردگی می‌شوند (ذخیره‌های وب ناقص هستند و برخی وب‌گاه‌ها در آرشیوها همزمان ظاهر نمی‌شوند). مانند سرخوردگی عملکرد جستجو: برای کارکنان، به ویژه، فقدان یک فهرست و این حقیقت که کتابداران هر کار نظام مند توصیفی را در مورد وب‌گاه‌ها انجام

ص: 299

نمی دهند برای آنان سخت کرده است که آرشیوهای وب را به عنوان مجموعه های میراث با ارزش مورد توجه قرار دهند (چقدر احتمال دارد یک کتابدار بتواند در مورد یک مجموعه ساخته شده توسط نرم افزار خزش گر و نه توسط خودش احساس مثبتی داشته باشد؟)

بنابراین چالش این است که کیفیت های متمایز آرشیوهای وب - به خصوص این حقیقت که آن ها تنها گواه واحد باقیمانده از تغییرات اساسی جامعه ما هستند که در 15 سال گذشته تجربه شده اند و گذارشان از آنالوگ به دیجیتال - را مشخص کرد، اما همچنین مشخص کردن الگوهای این مجموعه که در واقع آن را کمی شبیه به مجموعه های منظمی می کند که کتابداران یاد گرفته اند در دست بگیرند: مربوط می شود به شعار "دیجیتال متفاوت نیست" - و واژه را منتشر می کند.

گفته می شود که تجارت معمول آرشیو وب ساختن ابزار اولیه استفاده از ارتباط استاندارد، سازماندهی و راهبردهای بازاریابی است.

راهبردهای به حساب آوردن آرشیوهای وب به عنوان بخشی از کار روزانه کتابخانه

سازمان. مدیریت آرشیوهای وب نباید فقط به عنوان یک فعالیت فنی در نظام رسمی کتابخانه دیده شود بیش تر سازمان هایی که به آرشیو کردن وب به عنوان یک مورد فنی نگاه کرده اند، موضوعی تحت آی.تی. که نیاز به مهندس و رهبری توسعه نرم افزار دارد به سختی می توانند کتابداران را درگیر ارتقای آرشیوها کنند برای این که کار آرشیو کردن وب کاری معمول به حساب آید کتابخانه ملی فرانسه تصمیم گرفت فعالیت هایش را در قدیمی ترین واحد تولید، کتابخانه بخش واسپاری حقوقی به انجام برساند. وب گاه ها اکنون با منابع چاپی واسپاری حقوقی (کتاب ها و نشریات ادواری) در واحدی که سال 2008 به نام واحد «واسپاری حقوقی دیجیتال» (1) ایجاد شد، مدیریت می شوند. این واحد شامل گروهی پنج نفره می شود که فعالیت ها را با مشارکت متخصصان آی.تی. از یک طرف و متخصصان مجموعه از سویی دیگر اجرا می کنند این اجرا در بخشی بزرگ و قدیمی کمک بسیاری کرده است تا متخصصان و فعالیت های آرشیو کردن وب را در رسالت اصلی و تاریخ کتابخانه ملی فرانسه بگنجانند. به عبارت دیگر، این نوع سازماندهی به اتصال قدیم به جدید کمک کرده است و نمایش می دهد که راهبردهای نمونه برداری با مقیاس بزرگ برای حجم های عظیم داده های دیجیتال متولد شده خیلی متفاوت نیست از راهی که تاریخ سنت واسپاری حقوقی منابع چاپی فرانسوی را طی پنج سال گذشته شکل داده است. این نگرش بخشی از یک تلاش وسیع تر توسط کتابخانه ملی فرانسه است تا سازماندهی مهارت های کارکنان را با تغییر دیجیتال انطباق دهد. (2)

شبکه سازی. ساخت شبکه ای از متخصصان مجموعه موضوعی از سراسر کتابخانه مرحله سرنوشت ساز دیگری جهت تشویق مشارکت های فعال و پذیرش آرشیوهای وب به عنوان مجموعه ای

ص: 300

1- digital legal deposit

2- 'The Human Face of Digital Preservation: Organizational and Staff Challenges, and Initiatives at the Bibliothèque nationale de France', Emmanuelle Bermès, Louise Faudet, iPres 2009, [http://www.escholarship.org/uc/cdl_ipres09 [last accessed on 2010-06-15].

توسط کتابداران و مدیران بود. در سال 2005 این شبکه به عنوان گروهی متشکل از 20 پیشاهنگ شروع به کار کرد. تا سال 2010 حدود 80 کتابدار به گونه ای درگیر جمع آوری گزینشی بودند، به عنوان مثال، گزینش وبگاه کنترل کیفی و ارتقاء منابع و خدمات برای عموم نخست برای مدیران بخش مجموعه آگاهی ایجاد شد. در نتیجه، بیش تر بخش های موضوعی کتابداران موضوعی را برای طرح آرشیو کردن وب اختصاص دادند اکنون در هر حوزه اصلی (هنر ادبیات فلسفه علوم و غیره) از جمله موضوعاتی که مجموعه های نادر را در بر می گیرد (مانند نقشه ها یا موسیقی) یک کارگزار واسپاری حقوقی وب (1) وجود دارد، فردی که از آموزش خاص بهره مند بوده و دانش اساسی آرشیو کردن وب را کسب کرده است طوری که بتواند در حوزه تخصصی از آن استفاده کند برخی کتابداران موضوعی کتابخانه ملی فرانسه اکنون کتاب ها و وب گاه ها را به دست می آورند از سال 2008 به بعد تعداد کارکنان این شبکه کارگزاران مرحله واسپاری حقوقی وب که به خصوص به عنوان کتابداران مرجع در نظر گرفته شده بودند، افزایش یافته است. این جا هدف توسعه دانش و علاقه به آرشیوهای وب بود در هر کس که با عامه مردم خواه در اتاق های مطالعه یا به صورت آنلاین در ارتباط است.

ابزارهایی برای اشاعه داخلی. کتابخانه ملی فرانسه خوش شانس است که قادر است از مجموعه منابع مهمی برای ارتقای اطلاعات و مهارت ها به طور بین المللی سود برد: آموزش کارکنان، گفتگوهای داخلی ماهانه، کارکنان و کنفرانس هایی در وقت، ناهار یک مجله چاپی داخلی ماهانه و البته اینترنت همگی، قبل از شروع برنامه آرشیو کردن وب در دسترس بودند. چالش این بود که از تمامی این منابع ارتباطی استفاده و آن ها را به مسیری هوشمند هدایت کند تا کارکنان در مورد آرشیوهای وب در زمان و با روش مناسب بشنوند به عنوان مثال کتابخانه ملی فرانسه درست پایان مبارزات ریاست جمهوری سال 2007 یک نمایش اساسی داخلی از مجموعه وب گاه های منتخب برگزار کرد. همچنین مقالات مجله داخلی یا مقالاتی که در اینترنت کتابخانه ملی فرانسه بود یک طرح ارتباطی ای را دنبال کرد که منافع و علائق کلی کتابخانه را به حساب می آورد به عنوان مثال زمانی که کتابخانه مجموعه فعالیت هایی را که در راستای بالا بردن آگاهی نسبت به توسعه پایدار بود عملی کرد (و کارکنان را تشویق کرد که رفتاری مسئولیت پذیر در این حوزه به عهده بگیرند) گروه آرشیو کردن وب ارتباطات داخلی اش را روی مجموعه های مرتبط با توسعه پایدار نیز افزایش داد در یک کلام طرح ارتباط داخلی جهت ارتقاء آرشیوها هرگز به عنوان یک فرایند مستقل دیده نشد. ما تمامی فرصت های ممکن را به کار گرفتیم تا آرشیوهای وب رل در برنامه کلی کار کتابخانه بگنجانیم بیش از آن که تلاش کنیم این برنامه کار را با تحمیلی دیدن آن بر هم زنیم.

راهبردهایی برای رسیدن به کاربران نهایی

بحث ها: «آیا اصلاً آرشیوهای وب استفاده می شوند؟» «چرا منابع را صرف این داده های بیهوده کنیم در حالی که ما هیچ تضمینی نداریم که آیا این نوع ماده مورد علاقه عموم قرار خواهد گرفت؟» «چرا تلاش های مان را روی فراهم آوری یا دیجیتالی کردن موادی که می دانیم بدون شک میراثی برای نسل های

ص: 301

آینده، هستند متمرکز نمی کنیم؟» - این ها سؤالات معمولی است که خواه مربوط به مدیریت باشند یا مربوط به رسانه ها بیشتر اوقات نیاز به پاسخ دارند. از جهتی با نگاهی به گذشته تاریخ کتابخانه، چنین سؤال هایی جدید نیست در طول، زمان هر رسانه جدیدی سؤالات مشابهی برانگیخته است. دلیل آن این است که این امر زمان می برد - زمانی برای چیزهایی که از عموم فضای کسب و کار رایج ناپدید می شوند - قبل از این که پژوهشگران یا افراد متوجه شوند آن ها را از دست می دهند آنان به این مستندات نیاز دارند تا تاریخ جامعه را تشریح کنند.

روی هم رفته، موقعیت ویژه ای برای یک سازمان میراثی است تا برنامه آرشیو کردن وب خود را اجرا کند: گرایش به این وجود دارد که بروندهای ملموس را با مفید و با ارزش نشان دادن آن ها نمایش دهد در حالی که همزمان کاربرد نمی تواند بلافاصله توسعه یابد هنوز برای این کار خیلی زود است زیرا بیش تر کاربران آرشیوهای وب هنوز دارند متولد می شوند راهبردهای کتابخانه ملی فرانسه جهت کنترل این موقعیت دو جانبه است از یک سو با ایجاد ابزارها و آمارهایی که کتابخانه را قادر می سازد توسعه کاربرد را به روشی مشابه همان گونه که برای مجموعه های عادی انتظار می رفت نمایش دهد (دیجیتالی یا غیر دیجیتالی؛ به بخش بالا مراجعه کنید)، از سوی دیگر توجه عموم را به بحث در مورد قابلیت رؤیت از خارج از کتابخانه جلب کند.

ایجاد یک بحث عمومی. گر چه تعداد زیادی از مردم اکنون از آرشیوهای وب استفاده نمی کنند، در واقع بیشتر آنان علاقه به این مسأله را زمانی که از آن ها در این مورد سؤال شد، نشان داده اند. دلیل این امر این است که اینترنت زندگی های شخصی و عمومی افراد را تحت تأثیر قرار داده است بنابراین هر کس چیزی دارد که درباره آن بگوید. نخستین واکنش به طور معمول این است: «من هرگز در مورد آن فکر نکردم اما اکنون که شما به من می گوئید تصور می کنم مهم است این حافظه ماست و هر روز از یاد می رود» ارتباطات کتابخانه ملی فرانسه این مردم را برای توسعه حمایت عمومی از برنامه آرشیو کردن وب کتابخانه هدف قرار می دهد - آنان که ممکن است امروز برای دسترسی به مجموعه ها سروکله شان در کتابخانه پیدا نشود، اما این آگاهی را به دست آوردند که کاری است که باید انجام شود.

زنجیره کنفرانس حافظه های وب (1) به عنوان یک گردهمایی عمومی جهت فراهم آوردن این نوع حمایت و قابلیت رؤیت طراحی و در مارچ 2010 اجرا شد هر کنفرانسی یک نصف روز طول می کشد و سه نوع شرکت کننده را که در ترکیب مخاطبان منعکس می شوند نیز گردهم می آورد: پژوهشگران، وب ناشران و مسئولان همه کسانی که از آن ها درخواست می شود به موضوع مشابه بپردازند به عنوان مثال، عمل گرایی و سیاست وب خاطرات، وب حمایت از داده های شخصی در مقابل وسعت فضای عمومی، و غیره. کتابخانه ملی فرانسه نزدیک به 100 شرکت کننده را برای هر یک از این رویدادها گردهم آورده است و پوشش رسانه ای مناسبی را (بیشتر بر روی وب) برای دو نشست نخست دریافت کرده است. هدف این است که آرشیو کردن وب یک موضوع بحث عمومی در خارج از کتابخانه اما به همراه کتابخانه شود. این شکل ارتقا (که باز خورد عالی ای را برای ساخت مجموعه و خدمات به پژوهشگران نیز به دنبال دارد)

ص: 302

فعالیت های ارتقاء سازمان یافته مستقیم تر را برای کاربران نهایی امروز کتابخانه ملی فرانسه تکمیل می کند.

بیشتر کاربران آرشیوهای وب هنوز وجود ندارند آنان به احتمال با نسل های آینده ای می آیند که به دنیا آمده اند و با وب به دنیا خواهند آمد و آنان که از آن استفاده می کنند یا آن را به عنوان ابزار اصلی اطلاعاتی و ارتباطی شان استفاده خواهند کرد برای چنین مجموعه ای باید بپذیریم که فقط زمان است که به ما خواهد گفت آیا ما انتخاب درست را انجام می دهیم مدیریت دسترسی و گزینش نیز به معنی مدیریت خطرات است. ما مجبوریم آن چه که امروز غیر منتظره و ناخواسته است اما فردا ممکن است مورد علاقه، باشد بگیریم و ارتقا دهیم. همچنین مجبوریم کاربران و کاربرد فردا را در نظر بگیریم. این مسأله جدید نیست: چالشی است که مؤسسه میراثی قدیمی با آن آشناست. در این میان مؤسسات هنوز به راهبردهای میان مدت برای ایجاد اطمینان و حس مالکیت جمعی نسبت به این مجموعه های جدید در میان کارکنان و سرمایه داران شان نیاز دارند سازماندهی اشاعه داخلی ارتباطات و تلاش های آموزشی همگی در کل، کلید توسعه جوامع جدید آماده انطباق با مجموعه دیجیتال هستند.

منابع

1. Aubry, Sara (2008): 'Les archives de l'Internet: un nouveau service de la BnF', *Documentaliste Sciences de l'Information*, 45(4), pp. 12-13.
2. Bermès, Emmanuelle and Gildas Illien (2009): 'Metrics and Strategies for Web Heritage Management and Preservation', IFLA, <http://www.ifa.org/files/hq/papers/ifa75/92-bermes-en.pdf>
3. Bermès, Emmanuelle and Louise Faudet (2009): 'The Human Face of Digital Preservation: Organizational and Staff Challenges, and Initiatives at the Bibliothèque nationale de France', *iPres*, <http://www.escholarship.org/uc/cdl-ipres09>
4. Illien, Gildas (2008a): 'L'archivage d'Internet, un défi pour les décideurs et les bibliothécaires: scénarios d'organisation et d'évaluation, l'expérience du consortium IIPC et de la BnF', IFLA, <http://archive.ifa.org/IV/ifa74/papers/107-Illien-fr.pdf>
5. Illien, Gildas (2008b): 'Re-Inventing Collection Development Policy in the Age of Web Archiving: the Experience of the BnF', LIBER Annual Conference, <http://www.ku.edu.tr/ku/images/LIBER/LIBER-ILLIEN-2008.ppt>
6. Lasfargues, France, Clément Oury and Bert Wendland (2008): 'Legal deposit of the French Web: harvesting strategies for a national domain', IAWW, <http://iwaw.net/08/IWAW2008-Lasfargues.pdf>

WCT نوعی منبع باز برای مدیریت آرشیو وب گزینشی است که به عنوان پروژه ای مشترک بین کتابخانه ملی زلاندنو و کتابخانه بریتانیا توسعه یافته است. از ژانویه 2007 WCT در کتابخانه ملی زلاندنو به طور روزانه استفاده می شود. این مقاله نخستین سال آرشیو وب گزینشی ما را با نرم افزار جدید توصیف می کند. کتابخانه ملی زلاندنو، فواید WCT را توسعه داده و قصد دارد برنامه گزینشی گردآوری با WCT را برای آینده ای قابل پیش بینی ادامه دهد.

یک سال آرشیو وب گزینشی با (1) WCT در کتابخانه ملی زلاندنو

نوشته: گوردون پنیر، (2) سوزانا جو (3)، وانیتا لا لا (4)، گیلیان لی (5)

ترجمه: احترام السادات کیان مهر (6)

مقدمه

WCT نرم افزاری است که منابع برخط گزینش شده، گردآوری شده، و ارزیابی کیفی شده را که توسط کاربران گروهی در محیط کتابخانه ای به کار گرفته می شود پشتیبانی میکند از این نرم افزار، در مواقعی برای گردآوری وب گزینشی استفاده می شود که یک متخصص موضوعی قسمت هایی از یک وبگاه یا کل آن را معمولا در ارتباط با یک ناحیه موضوعی تمرکز یافته یا یک رویداد مهم شناسایی کرده باشد.

نرم افزار، به عنوان پروژه ای گروهی بین کتابخانه ملی زلاندنو و کتابخانه بریتانیا توسعه داده شده است و زیر نظر کنسرسیوم فراهم آوری اینترنت بین المللی مدیریت می شود.

WCT، نرم افزار منبع باز است و از طریق وبگاه <http://webcurator.sf.net> برای استفاده جامعه آرشیو وب بین المللی آزادانه قابل دسترس است.

کتابخانه ملی زلاندنو (از این پس کتابخانه) از ژانویه 2007 از WCT به عنوان پایه برنامه آرشیو وب

ص: 305

Web Curator Tool -1

Gardon Panynter -2

Susanna Joe -3

Vanita Lala -4

Gillian Lee -5

-6 کارشناس ارشد سازمان اسناد و کتابخانه ملی

گزینشی استفاده می کند.

طی سال اول ویرایش جدید نرم افزار توسعه داده شد و به طور شگرفی، افزایش و بهبود کیفیت فعالیت های گردآوری و نیز گردآوری منابع و بی به طور خودکار را میسر ساخته است.

این مقاله تجربه ما را در استفاده از WCT در محیط کار تعریف می کند

بخش بعدی، مقاله پیش زمینه فعالیت های بهره برداری وب کتابخانه و WCT را در اختیار می گذارد بخش های زیر تجربه ما را با نرم افزار جمع بندی می کند و سرانجام با توصیف یک رویداد اجرا شده گردآوری با نرم افزار خاتمه می یابد.

2. آرشیو وب گزینشی در کتابخانه ملی زلاندنو

2-1. انگیزه

کتابخانه ملی زلاندنو حکمی قانونی و مسئولیتی اجتماعی برای محافظت تاریخ فرهنگی و اجتماعی زلاندنو به شکل، کتاب، روزنامه عکس وبگاه بلاگ، یوتیوب، و ویدئو بر عهده دارد.

علاوه بر آن میراث مستند کتابخانه ملی زلاندنو فقط به صورت برخط قابل دسترس است از نظر این محتوای دسترس پذیر با ارزش است اما ناپایداری فقدان مالکیت شفاف و صحیح و طبیعت پویای آن چالش های مهمی برای هر مؤسسه ای است که برای اندوختن و نگهداری آن ها تلاش می کند.

وب گاه ها، WCT برای حل همین مشکلات ارتقا یافت به این ترتیب که به برخی مؤسسه ها اجازه می داد هر نوع سند برخط شامل صفحه های وب، وب گاه ها و وب نوشت ها و سایر اشکال جاری، شامل صفحه های HTML، تصاویر PDF و مدارک word مانند محتوای چند رسانه ای نظیر فایل های دیداری و شنیداری را دریافت کنند این انواع با مراقبت های، ممکن طوری هماهنگ می شوند که یکپارچگی و اصالت شان حفظ شود

برخط منفعت عمومی از نگهداری و حفاظت دراز مدت میراث برخط زلاندنو غیر قابل محاسبه است. تاریخ اجتماعی برخط ما و اکثر تاریخ دولت و سازمان برای حفظ کردن قابل امکان خواهد بود و برای محققان تاریخ دانان و شهروندان آینده زلاندنو نیز امکان پذیر خواهد بود آن ها خواهند توانست به گذشته مستند مدارک دیجیتال ما در راهی مشابه که زلاندنویی ها تا امروز به واژه های چاپ شده که برای ما از نسل های گذشته باقی مانده است، نگاه کنند.

2-2. تاریخچه درو / هاروستینگ

کتابخانه ملی، زلاندنو از سال 1999 تا پایان سال 2006، برنامه عملی آرشیو وب گزینشی را اجرا کرده است، کتابخانه برای نگهداری منابع از نرم افزار کپی (1) استفاده کرده و منابع را بر اساس مارک (2) هدایت و پایگاه داده ها فراهم می کند. نرم افزار HTTrack کتابخانه را با یک پس افت منابع گردآوری شده که

HTTrack Website Copier Software -1
(MARC (Machine readable cataloging -2

قابلیت آرشیو شدن برای نگهداری دراز مدت را ندارد رها می کند. در حال حاضر، برای تبدیل این منابع به شکل مناسب قابل آرشیو کردن برنامه انتقال داده در حال اجراست.

2-3. نرم افزار گردآوری وب

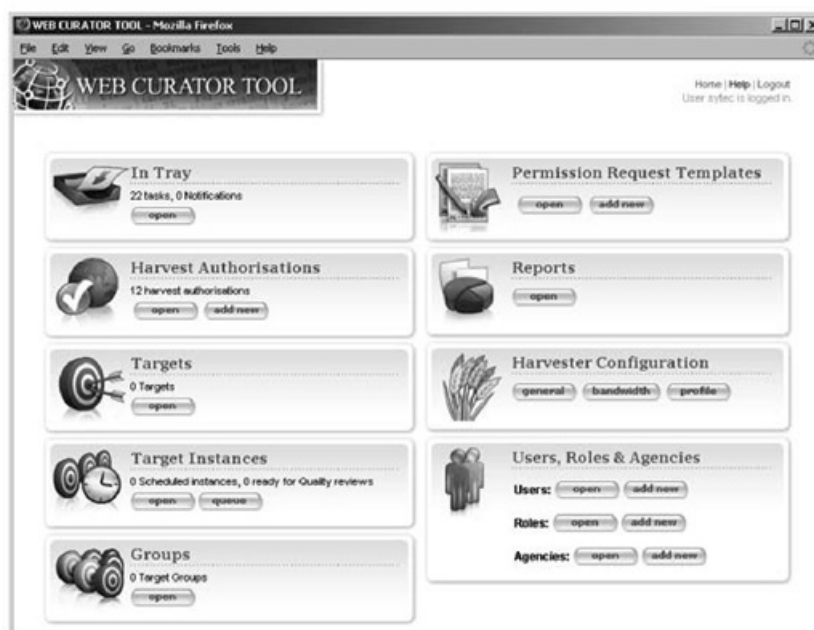
WCT، نوعی نمودار کاری گردآوری را پشتیبانی می کند که شامل وظایف تخصصی است: انتخاب منبع برخط، جست و جو، اجازه نگهداری، آن قابل دسترس بودن، قابل توصیف بودن، تعیین دامنه آن، و دسته بندی ها، نمودار کردن گردآوری وب یا یک سری از گردآوری های وب، اجرا کردن گردآوری ها، انجام دادن مرور کیفیت و تأیید یا تکذیب کردن منابع گردآوری شده و ذخیره منابع تأیید شده در یک مخزن دیجیتال یا آرشیو.

عکس

قابلیت آرشیو شدن، برای نگهداری دراز مدت را ندارد، رها می‌کند. در حال حاضر، برای تبدیل این منابع به شکل مناسب قابل آرشیو کردن برنامه انتقال داده در حال اجراست.

۲-۳. نرم افزار گردآوری وب

WCT، نوعی نمودار کاری گردآوری را پشتیبانی می‌کند که شامل وظایف تخصصی است: انتخاب منبع برخط، جست‌وجو، اجازه نگهداری آن، قابل دسترس بودن، قابل توصیف بودن، تعیین دامنه آن، و دسته‌بندی‌ها، نمودار کردن گردآوری وب یا یک سری از گردآوری‌های وب، اجرا کردن گردآوری‌ها، انجام دادن مرور کیفیت و تأیید یا تکذیب کردن منابع گردآوری شده و ذخیره منابع تأیید شده در یک مخزن دیجیتال یا آرشیو.



شکل ۱. فهرست اصلی نرم‌افزار را نشان می‌دهد.

شکل ۱. فهرست WCT

بیشتر فعالیت‌های آرشیو وب به‌طور عمده‌ای بر تخصص‌های فنی اپراتورهای گردآوری تکیه دارد. از سوی دیگر، WCT گردآوری را مسئولیت کاربران و متخصصان موضوعی (ترجیحاً از مهندسان و مدیران سیستم) از طریق آسان و روان نمودن، به‌طور خودکار، جزئیات فنی گردآوری وب می‌داند.

شکل ۱. فهرست اصلی نرم‌افزار را نشان می‌دهد.

شکل ۱. فهرست WCT

بیشتر فعالیت‌های آرشیو وب به‌طور عمده‌ای بر تخصص‌های فنی اپراتورهای گردآوری تکیه دارد. از سوی دیگر، WCT گردآوری را مسئولیت کاربران و متخصصان موضوعی (ترجیحاً از مهندسان و مدیران سیستم) از طریق آسان و روان نمودن، به‌طور خودکار، جزئیات فنی گردآوری وب می‌داند.

نرم افزار برای کار مطمئن و تأثیر گذار در محیط جایی که کارکنان پشتیبانی فنی می توانند آن را نگهداری کنند طراحی شده است.

WCT نرم افزار منبع باز است و به طور آزادانه از طریق وبگاه <http://webcurato.sf.net> تحت گواهینامه Apache public License قابل دسترس است.

وبگاه دسترسی کاربران به، راهنماها دست نامه ها فهرست های، پستی screenshots، سؤال ها، مستندات فنی و اداری کد منبع ردیابی اشتباهات و صفحه پروژه sourceforg را فراهم می کند.

(پینتر و میسون (1)، نرم افزار و توسعه اش را با جزئیات بیش تری در مقاله خود در کنفرانس لیانسا (2) در سال 2006 توصیف می کنند).

4.2. اعضا و منابع

در کتابخانه ملی زلاندنو، WCT نرم افزار اولیه و مسئولیت کتابدارن نشر الکترونیکی در کتابخانه «الکساندر تورن بل» (3) است. در سال 2007 معادل دو و نیم برابر انتخابگر، الکترونیکی به طور مستقیم از این نرم افزار استفاده کردند و نیز همه گزینش ها را مدیریت نموده و گردآوری و مرور کیفی نمودند.

به هر حال نرم افزار با خط مشی های کتابخانه گردش کار ارتباطات راه دور و پشتیبانی سرویس ها یکپارچه و جامع است و بر گروه گسترده تر و بیش تری از اعضا تأثیر می گذارد.

برای مثال، سخت افزار و نرم افزار توسط سرویس های فنی حفظ می شوند و توسط میز پشتیبانی مدیریت می شوند، رده بندی توسط سرویس های محتوا هدایت می شود و آرشیو دیجیتال توسط کتابخانه ملی دیجیتال حفظ و نگهداری می شود.

WCT برای یکپارچه شدن با سیستم های کاری، موجود به طور محکم و استوار - و تا حد ممکن - طراحی شده است.

برای استقرار وضعیت از سخت افزار استاندارد کتابخانه (سرورهای سان اسپارک (4) سیستم های عملیاتی (سولاریس (5))، پایگاه اطلاعاتی (اوراکل (6))، وب سرویس ها (Apache HTTP Server and Tomcat)، و خدمات ثبت کاربر (راهنمای الکترونیکی نوول (7) استفاده می شود.

سیستم تولید بین دو سرور مستقر شده، است یکی برای مدل کور (8) و دیگری برای نرم افزار گردآوری (هم آرایه کردن برای گردآوری هشت منبع و بی همزمان) و اشتراک پایگاه اطلاعاتی موجود و سرورهای فایل با سیستم های دیگر کتابخانه سیستم آزمون جداگانه با ترتیب و وضع مشابه حفظ می شود.

ص: 308

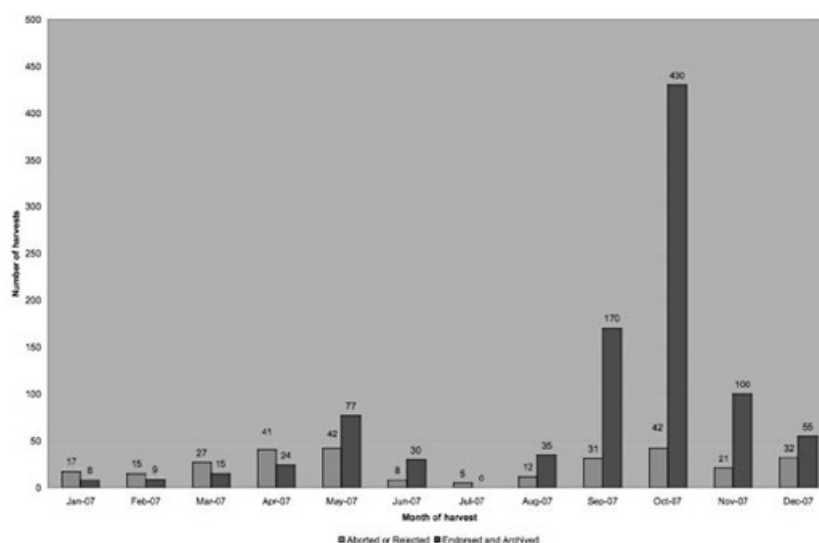
Sun SAARC Servers	-4
(Operating System (Solaris	-5
Oracle	-6
Novell e Directory	-7
Core module	-8

نمودار 2 تعداد گردآوری های خودکار انجام شده توسط کتابخانه در سال 2007 را نشان می دهد. برای هر ماه، ستون ارغوانی رنگ شماره نتایج را نشان می دهد که در آرشیو تأیید نشده است. این موارد یا متوقف شده‌اند (قبل از آن که بتوانند کامل شوند) یا طی مراحل کنترل کیفیت به دلیل آن که برخی جنبه های اساسی وبگاه را کسب نکرده اند رها شده اند ستونهای قرمز تعداد نتایج موفقیت آمیزی که در آرشیو کتابخانه دیجیتال آرشیو شده اند نشان می دهد. در مجموع 1249 گردآوری منابع خودکار توسط سیستم محصول ما در سال 2007 کمک شده‌اند. 953 (76 درصد) به طور موفقیت آمیزی گردآوری شده، کنترل کیفی شده و نیز آرشیو شده است.

عکس

۲-۵. سطح گردآوری / درو

نمودار ۲، تعداد گردآوری‌های خودکار انجام شده توسط کتابخانه در سال ۲۰۰۷ را نشان می‌دهد. برای هر ماه، ستون ارغوانی رنگ شماره نتایج را نشان می‌دهد که در آرشیو تأیید نشده است. این موارد یا متوقف شده‌اند (قبل از آنکه بتوانند کامل شوند) یا طی مراحل کنترل کیفیت، به دلیل آنکه برخی جنبه‌های اساسی وبگاه را کسب نکرده‌اند رها شده‌اند. ستون‌های قرمز، تعداد نتایج موفقیت‌آمیزی که در آرشیو کتابخانه دیجیتال آرشیو شده‌اند نشان می‌دهد. در مجموع ۱۲۴۹ گردآوری منابع خودکار توسط سیستم محصول ما در سال ۲۰۰۷ کمک شده‌اند. ۹۵۳ (۷۶ درصد) به‌طور موفقیت‌آمیزی گردآوری شده، کنترل کیفی شده و نیز آرشیو شده است.



نمودار ۲. سطح گردآوری گزینش شده در ۲۰۰۷

باقیمانندگان ۶۹ (۶ درصد) طی گردآوری متوقف شده و ۲۲۴ (۱۸ درصد) رد شده‌اند، زیرا از مرحله کنترل کیفی عبور نکرده‌اند.

سه فاز گردآوری و برداشت منابع وبی به‌طور خودکار در نمودار ۲ قابل رؤیت هستند. نخستین آنها از ژانویه تا می است، یعنی زمانی که انتخابگر الکترونیکی در فعالیتهای گردآوری از ویرایش ۱.۱ WCT استفاده کرده است. گردآوری در این مرحله کم است و در ماه می ۷۵ گردآوری موفق داشته است، یعنی زمانی که گروه در حال تجربه و گردش کار و خروجی‌های مختلف - حتی آنهایی که هنگام کار با نرم‌افزار اولیه کشف نشده بود - بودند.

نمودار 2. سطح گردآوری گزینش شده در 2007

باقیمانندگان 69 (6 درصد) طی گردآوری متوقف شده و 224 (18 درصد) رد شده‌اند، زیرا از مرحله کنترل کیفی عبور نکرده‌اند.

سه فاز گردآوری و برداشت منابع وبی به‌طور خودکار در نمودار 2 قابل رؤیت هستند. نخستین آن‌ها از ژانویه تا می است یعنی زمانی که انتخابگر الکترونیکی در فعالیتهای گردآوری از ویرایش 1.1 WCT استفاده کرده است. گردآوری در این مرحله کم است و در ماه می 75 گردآوری موفق داشته است، یعنی زمانی که گروه در حال تجربه و گردش کار و خروجی‌های مختلف - حتی آن‌هایی که هنگام

کار با نرم افزار اولیه کشف نشده بود - بودند.

ص: 309

از ژوئن تا اواسط سپتامبر، انتخابگران الکترونیکی آزمون ویرایش WCT 1/2 را تقاضا کردند. بنابراین، مقررات فعالیت های منظم گردآوری حتی بیشتر از قبل قطع می شود.

فاز نهایی گردآوری در اواسط سپتامبر شروع می شود و افزایش شگرفی در فعالیت گردآوری را نشان می دهد. عوامل متعددی در این افزایش دخیل هستند که شامل انتخاب ویرایش WCT 2/1 در سیستم، تولید افزایش 50 برابری ظرفیت اینترنت کتابخانه (از 2 مگابایت در ثانیه به 100 مگابایت در ثانیه)، شفافیت بیشتری خروجی های گردش کار پیرامون گردآوری کردن و فهرست کردن و آغاز دورویداد بر اساس گردآوری (در زیر توضیح داده شده).

در نتیجه، تقریباً گردآوری در نیمی از سال فقط در ماه (تعداد وب گاه ها) اکتبر انجام می شود.

نکته جالب در این مرحله این است که اگر چه تعداد گردآوری های موفق به طور چشمگیری افزایش یافته است تعداد «ناتمام» ها و «مردود» ها افزایش نیافته است و عددی که تأیید یا رد می شود، محسوب نمی گردد.

گردآوری با WCT ویرایش 1.1.

زمانی که انتخاب گر الکترونیکی از ویرایش 1.1.1 WCT از ژانویه تا می استفاده کرد، از تجربه با این نرم افزار سود بردند و به فعالیت های گردآوری و نیز مشکلات مختلفی که از طریق ویرایش 2.1 و کامل شدن آن بود پی بردند

این بخش خلاصه ای از این تجربه است و توصیف بیشتر جزئیات از نویسندگان بر اساس درخواست در دسترس می باشد.

1.3. تجربه اولیه

مهم ترین تغییر ایجاد شده توسط WCT این است که نرم افزار منابع گردآوری شده را در فرمت فایل ARC ذخیره می کند که از لحاظ فراهم آوری سودمند است اما برای فعالیت های کنترل کیفیت.

و این موضوع راه اساسی جدیدی را برای انتخابگران الکترونیکی برای رسیدن و نزدیک شدن به تجدید نظر کیفی را به همراه دارد.

با این شیوه جدید فعل و انفعال، مشکل می توان گفت که آیا خروجی های کیفی، توسط نرم افزارهای گردآوری کننده ایجاد شده یا توسط کمبودهای نرم افزارهای کنترل کیفی

برای مثال اگر یک بازیکن کننده یک صفحه وب را در نرم افزار تورق باز کند و نشان داده شود که یک صفحه وب (تصویر) گم شده است انتخابگر الکترونیکی ابتدا باید تعیین کند که آیا به طور موفق، تصویر از منابع و بی گردآوری شده یا آن را به طور موفقیت آمیزی گردآوری کرده، اما مرورگر نمی تواند آن را نشان دهد خیلی زود مشخص شد که برای پشتیبانی اعضا در وجود این تفاوت ها تمرین های اضافی ضروری بوده است

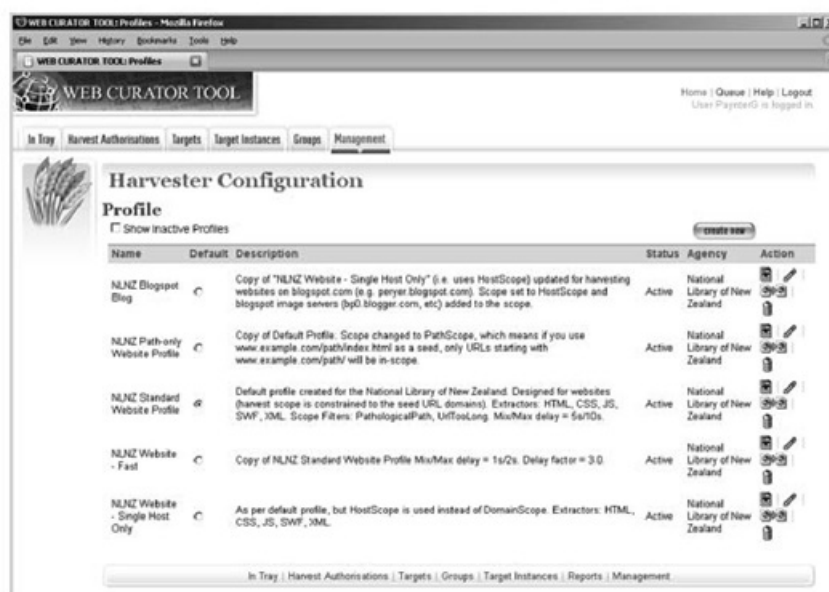
زمانی که کتابداران تجربه مهمی از مرور نتایج گردآوری نرم افزار HTTrack را کسب کردند، نمی دانستند که چگونه از عدم موفقیت گردآوری منابع وبی در WCT گره بکشایند.

تمام مشکلات جدید گردآوری در نسخه جدید ضبط و مشاهده شد که همان نوع اشتباه ها دوباره رخ داده است. بنابراین بخشی از تمرین حفظ و یک روند بازشناختی نیز توسعه و به WCT به طور دستی اضافه شد.

مشکل دیگر که در ماه های نخست نمایان شد این بود که توسعه سودمند پروفایل های گردآوری شده، به خصوص در مورد بلاگ ها زمان بر است توسعه شکل 3. مجموعه ای از پروفایل های گردآوری شده را نشان می دهد که سرانجام به وجود آمده اند.

عکس

زمانی که کتابداران تجربه مهمی از مرور نتایج گردآوری نرم افزار HTTrack را کسب کردند، نمی دانستند که چگونه از عدم موفقیت گردآوری منابع وبی در WCT گره بگشایند. تمام مشکلات جدید گردآوری در نسخه جدید ضبط و مشاهده شد که همان نوع اشتباهها دوباره رخ داده است. بنابراین، بخشی از تمرین حفظ و یک روند بازشناختی نیز توسعه و به WCT به طور دستی اضافه شد. مشکل دیگر که در ماههای نخست نمایان شد این بود که توسعه سودمند پروفایل های گردآوری شده، به خصوص در مورد بلاگها، زمان بر است. توسعه شکل ۳. مجموعه ای از پروفایل های گردآوری شده را نشان می دهد که سرانجام به وجود آمده اند.



شکل ۳. پروفایل های گردآوری کننده در حال استفاده در کتابخانه در دسامبر ۲۰۰۷

۳.۲. مشکلات گردآوری

مشکلات زیر در ویرایش ۱.۱ وجود داشته و در ویرایش ۲.۱ بر طرف شده است.

کمبودهای تورق نرم افزار

برخی اوقات محتوای ظاهر شده از دست می رود به دلیل اینکه نرم افزار تورق (مرور) نمی تواند آن را به طور صحیحی منتقل کند.

شکل 3. پروفایل های گردآوری کننده در حال استفاده در کتابخانه در دسامبر 2007

3.2. مشکلات گردآوری

مشکلات زیر در ویرایش 1.1 وجود داشته و در ویرایش 2.1 بر طرف شده است

کمبودهای تورق نرم افزار

برخی اوقات محتوای ظاهر شده از دست می‌رود به دلیل اینکه نرم افزار تورق (مرور) نمی‌تواند آن را به طور صحیحی منتقل کند.

ص: 311

برخی نسخه های کاربر برای انتخاب گران الکترونیکی، سیر کردن سریع در پیرامون نرم افزار را مشکل می ساخت.

اشتباه «در مرحله ایست درگیر شدن»

مشکل دوباره به وجود آمده در ویرایش جدید WCT این بود که گردآوری باید همان گونه که انتظار می رفت با موفقیت به مرحله پایانی می رسید ولی سیستم به وضعیت «گردآوری شده» انتقال نمی یافت و در عوض در وضعیت توقف باقی می ماند مشکلاتی که به این شرایط هدایت می شود نتایجی از نسخه های مورد نظر که در بالا آمده و کامل شده است.

وب گاه های بزرگ

کتابخانه بریتانیا به این موضوع پی برد که نرم افزار از عهده گردآوری های بزرگ و گسترده بر نمی آید و مشکلاتی با نرم افزار گردآوری همزمان وجود دارد.

3.3. خلاصه تجربه با ویرایش 1.1

با وجود برخی ای محدودیت ها ویرایش اول WCT بهبودی وسیع و بزرگی بر نرم افزار HTTrack قبلی بود و آن مرحله مهم و با اهمیتی بود وقتی که ما اولین وب گاه های گردآوری شده را به آرشیو دیجیتال ارائه دادیم

طی سال اول WCT 1250 منابع وبی را به طور خودکار گردآوری کرد و تقریباً 950 مورد دوباره، بررسی تأیید و به آرشیو اضافه شد.

گردآوری با WCT با ویرایش 1.2

ویرایش WCT 1.2.6، در اواسط سپتامبر 2007 به محصول گسترش داده شد و از آن در ماه های بعدی در دو رویداد گردآوری استفاده شد و سپس در پایان اکتبر با ویرایش 1.2.7 جابه جا شد.

4.1. نرم افزار گردآوری تاریخ

یکی از نواحی ای که در آن پیشرفت زیادی به وجود آمد کنترل کیفیت است. افزایش در این ناحیه به طور شگرفی تأثیر کنترل کیفیت را ثابت کرد و باعث پیشرفت بیشتر، سریع تر، ماهرانه تر شد.

نرم افزار گردآوری تاریخ در کنترل کیفی ثابت کرده است که بسیار مفیدتر از تقویت کنترل کیفی است. این نرم افزار همه منابع گردآوری شده را که دارای یک هدف ویژه است با چکیده اطلاعات، نظیر تاریخ شروع تاریخ انتقال URL گردآوری شده و افتاده، زمان صرف شده و وضعیت جاری را فهرست می کند.

این امر بسیار مفید است به طوری که از اطلاعاتی که در روند کنترل کیفیت نیاز بود محکم تر است.

2.4. نرم افزار تورق (مرورگر)

هنوز یک تغییر ساده و تأثیر گذار برای مرور میانجی، بود افزایش پیشنهادها برای دریافت دیداری یک سایت گردآوری شده از سه راه مختلف صورت می گیرد: دریافت دیداری مرحله جاری گردآوری شده در نرم افزار مرورگر دریافت دیداری سایت، جاری یا دریافت دیداری ویرایش آرشیو شده قبلی (در آرشیو اینترنت یا آرشیو مکانی).

این مراحل در Tab دیگر با تورق ویندوز باز می شود و اجازه می دهد که مرورگر یک کپی گردآوری شده را با ویرایش دیگر سایت مقایسه کند.

این پیشرفت بزرگ باعث شد که نرم افزار تورق به ما اجازه دهد که به طور موفقیت آمیزی مرور و تورق سایت های دیگر را - که قبلاً قابل دیدن نبودند - مشاهده کنیم به هر حال، چند خروجی حل نشده باقی می ماند.

وب گاه های تحویل داده شده که از جاوا اسکریپت استفاده می کند باعث مشکلاتی می شود، به ویژه دوباره عناصر وظیفه ای مانند منوی پایین - بالا که می تواند سایت های گردآوری شده را برای سیر کردن مشکل سازد.

گردآوری سبک برگه ها هنوز می تواند مشکل باشد، گرچه نسخه های جمع آوری شده بیش از نسخه های تکراری است نسخه های برجسته، دیگر شامل مشکلات تکرار شده URL ها با فضاها و جمع آوری تصاویر پیش زمینه جاسازی شده مطمئن.

3.4. نرم افزار هرس درخت تصمیم

نرم افزار هرس درخت تصمیم سیستم روز آمد شد اما کتابخانه نرم افزار هرس را در وب گاه های قدیمی تر استفاده نمی کند. بنابراین از این تغییرات استفاده زیادی به عمل نمی آید. به هر حال، ما با نرم افزار تجربه کرده ایم و مدتی می تواند خیلی آهسته برای سایت های بزرگ استفاده شود وظیفه جدید نرم افزار هرس ثابت کرد که می تواند برای انتخاب کردن و نگاه کردن فایل های دل خواه از منابع گردآوری شده بسیار مفید باشد

4.4. وب گاه های بزرگتر جمع آوری می شوند

در ماه های آخر سال ما متوجه شدیم که فعالیت های گردآوری ما شامل سایت های بزرگتر و بزرگ تر هستند. بزرگ ترین گردآوری، کتابخانه به عنوان قسمتی از برنامه گردآوری در اندازه 21 گیگا بایت کامل و کنترل شد گرچه پس از کنترل کیفیت رد شده بود.

وب گاه های متعدد 10 گیگا بایتی به طور موفقیت آمیزی کنترل آرشیو و گردآوری شدند زمانی که اندازه وب گاه ها رشد کرد، انتخاب گران الکترونیکی تکیه بیش تری بر پروفایل های عادی و پروفایل های

برتر کردند به ویژه فیلتر ها که به انتخاب گران الکترونیکی اجازه داد که گردآوری کننده را از جمع آوری قسمت های ویژه وبگاه متوقف کنند.

5.4. گسترش ذخیره دیجیتال با ارزش

نتیجه افزایش تعداد و اندازه وب گاه های گردآوری شده این بود که WCT به ذخیره دیجیتال با ارزش بیش تری نیاز داشت (فضای موقت برای ذخیره گردآوری ها زمانی که کنترل کیفی شده اند و قبل از این که آرشیو شوند).

در اکتبر و نوامبر این مشکلی ویژه و خاص بود زیرا تعداد وب گاه های انتخاب شده گردآوری شده و تأیید شده ناگهان افزایش یافت و موقتاً به طور طولانی تری فهرست شد و نوعی انبار منبع به وجود آمد. در نتیجه فضای اختصاص داده شده برای ذخیره منابع دیجیتال با ارزش که زودتر ظاهر شده از دست رفت و ناگهان به نظر رسید که بسیار کم ذخیره شده است در نتیجه یک دیسک بزرگ تر و جدید برای ذخیره با ارزش منابع دیجیتال بود و ابزارهای گزارش گیری به کتابداران اجازه می دهد که از دیسک به طور موقت استفاده کنند و نمایش دهند.

این وضعیت، همچنین نسخه ای از آرشیو موقتی دیجیتال را فراهم کرد که در آن فضای دیسک با ذخیره ارزشمند دیجیتال مشترک شده مطمئناً تفصیل کنترل های نهایی فایده ای نداشته است.

6.4. ارتباط

ارتباط خوب بین اعضای کتابخانه های مختلف از طریق نرم افزار تحت تأثیر واقع می شود، اما شکاف های مختلفی در ارتباط از زمانی که WCT از ابتدا استفاده می شد در محصول شناسایی شد.

گرچه در ویرایش 1.2 پیشرفت هایی حاصل شد خیلی از قلم افتادگی ها، خارج از نرم افزار از طریق گزارش های دوره ای از طریق پایگاه اطلاعاتی WCT نشانی شده اند (مدرک نرم افزار، شامل یک فرهنگ لغت داده کامل است).

برای مثال، گزارش هفتگی حاصل شده از گردآوری، منابع به فهرست نویس ها برای اعلام آن ها به وبگاه هایی که نیاز است فرستاده شد و تعداد و اندازه وب های گردآوری شده جاداده شده توسط انتشارات الکترونیکی کتابداران برای سرویس های فنی اعضا پیش بینی شد. زمانی که این ها در یک شکل کامل گردآوری شدند ما امیدواریم که به آن ها یک سری گزارش های ساختار WCT را اضافه کنیم.

7.4. فهرست گردش کار و دسترسی

زمانی ممکن است از WCT برای توصیف وب گاه ها استفاده شود و آن خط مشی کتابخانه است برای توصیف مجموعه کامل کتابخانه در فهرستش و فراهم کردن یک پیوند رکورد فهرست کتابخانه به آیم های دیجیتال نگهداری شده در محزن دیجیتال

(به هر حال انبارش های دیجیتال، پیچیده نظیر سریال ها و وب گاه ها به طور عام قابل دسترس نخواهند)

بود تا آرشیو میراث دیجیتال ملی با مخزن دیجیتال موقت ما در سال 2008 جایگزین شود).

یک مانع و اشکال در جست و جو کردن فهرست کتابخانه برای وب گاه هاست و ناتوانی جمع آوری کامل رویدادها از هر وبگاه که به تنهایی فهرست شده است

نشانی این نسخه کتابخانه بریتانیا توسعه یک نرم افزار میانجی وب را طراحی می کند که دسترسی به منابع گردآوری شده را بر اساس موضوع و رویداد به عنوان یک هدف اضافی برای کاربرانی که می خواهند به طور ویژه برای وب گاه ها جست و جو کنند فراهم می کند ما پیشرفت کتابخانه بریتانیا را با علاقه بسیار زیادی دنبال خواهیم کرد.

5. گردآوری رویداد انتخابات هیئت محلی

کتابخانه، با به کار بردن نرم افزار گردآوری HTTrack در سال 2007 مسئولیت گردآوری برخی رویدادهای عظیم ورزشی را بر عهده داشته است (جام آمریکا) (1) و انتخابات (انتخابات پارلمان ملی).

کتابخانه دو رویداد گردآوری را طراحی کرد که به طور اتفاقی با استقرار ویرایش WCT 1/2 همزمان بود. گزینش های هیئت دولتی محلی، رویدادی سه ساله است و هر قدرت محلی در زلاندنو به نگهداری و هدایت گزینش اعضایش احتیاج دارد کتابخانه برنامه گردآوری رویداد دوازده هفته ای را با تمرکز بر گزینش های هیئت محلی در سال 2007 و شروع شدن در سپتامبر را به عهده گرفت. گردآوری رویداد هیئت محلی اولین گردآوری رویداد عظیمی با استفاده از WCT بود و همچنین هنوز بزرگ ترین تلاش در کتابخانه بود با 238 وب گاه انتخاب شده

این سایت ها شامل وب گاه ها، بلاگ ها، وب گاه های شوراهای منطقه ای، شهری و ناحیه ای، اخبارهای سایت ها و سایت های عمومی یا بلاگ ها با محتوای مربوط و مناسب هم گزینشی یا دولت محلی با این حال طیف وسیعی از وب گاه ها انتخاب شده اند که تفسیرهای رسمی و غیر رسمی را با پراکندگی جغرافیایی وسیعی نمایش می دهند کلیه سایت های انتخاب شده در محدوده وضعیت قانونی زلاندنو بودند، بنابراین نیازی به جست و جوی دستور صریح و واضح و قانون گذاری برای به دست آوردن آن ها نداریم

سایت های انتخاب شده برای تعیین گردآوری، اولویت بندی شدند و در آن موقع یک جدول گردآوری برای پوشاندن کلیه تاریخ ها در مدت 16 هفته طراحی شده بود. جدول با تعداد زیادی سایت انتخاب شده و با قابلیت اداره خاتمه داده شد درگیری مقتضی اولیه برای پایان دادن منابع اعضا نیز وجود داشت. تعداد کمی از گردآوری ها یک فاصله زمان کافی بین هر گردآوری برای 2 عضو برای کامل کردن خزش گرها و روند تکرار کیفیت را مطمئن ساخت.

عامل تأثیر گذار دیگر بر برنامه گردآوری از گردش کاری هایی ریشه می گرفت که در حال حاضر از WCT استفاده می کنند در گردآوری های قبلی ایجاد نمودار گردآوری تکراری عظیمی از کنترل کیفیت امکان پذیر بود که می توان بعد از تکمیل همه گردآوری ها انجام داد.

ص: 315

به هر حال با تعداد زیادی از سایت های انتخاب شده در این دوره این رویکرد بر WCT هم غیر عملی و نشدنی است و هم سنگین و طاقت فرساست.

با گردآوری و سری زمانی، اهداف گروه ها با استفاده از سیستم های گروهی و وظیفه ای به وجود آمد. سایت های برگزیده مرتب شدند و 36 گروه با دو مقوله وبگاه، یا - در مورد سایت های شور- منطقه جغرافیایی تعیین شدند.

این گروه ها غیر رسمی بودند، زیرا استفاده اولیه گروه ها در این رویداد (دوره) - برای آسان تر شدن برنامه گردآوری صورت می گرفت و نه برای گروه بندی رسمی استفاده از گروه در مدیریت گردش کاری در تیم مفید بود و متناوب شدن گردآوری توسط گروه ها نیز مفید بود. همچنین نوعی تعادل فشار در بین سخت افزارهای گردآوری به وجود آورد.

اگر چه هدف، شروع گردآوری در بعد از ظهر بود، حتی با وجود هشت گردآورنده همزمان این کار ممکن نشد برخی گردآوری های برنامه ریزی شده چندین ساعت در نوبت می ماندند تا یک گردآورنده آزاد شود، اما خوشبختانه صف های طولانی برداشت موردی به وجود نیاورد و با کوچک شدن بسیاری از سایت ها صف های برداشت به سرعت ناپدید شد این سیستم با استفاده از هشت گردآورنده همزمان به خوبی با دوره های گردآوری طولانی و فشرده، مقابله کرد

مشکلات وقتی به وجود آمد که خزش های وبگاه شوراها در اندازه های بسیار بزرگتر با آرشیو دیجیتال وب گاه هایی که منتظر فهرست نویسی پیش آرشیوی بودند همزمان شدند. نتیجه یک تفسیر مشکل بحرانی در ظرفیت فضای دیسک بود. از آن جا که دیگر مناطق کتابخانه به طور مشترک از همین امکان ذخیره سازی استفاده می کردند عواقب بسیار جدی داشته است. در یک مقطع مجبور شدیم گردآوری را تازمانی که فضای دیسک کافی در دسترس باشد به تعویق بیندازیم.

با وجود پیچیدگی های اخیر، بیش تر وب گاه هایی جمع آوری شده بررسی کیفی و تأیید شده و با موفقیت آرشیو شده اند. در پایان برنامه برداشت، بیش از 600 وبگاه برداشته شده تأیید و یا برای این رویداد خاص آرشیو شده اند که موفقیت (دستاورد) قابل توجهی برای برنامه آرشیو وب کتابخانه است.

به طور کلی انتخابات هیئت، محلی آزمایش مهمی درباره گردآوری رویداد با استفاده از WCT بود. ما به سهولت قادر به برنامه ریزی و مدیریت این رویداد بزرگ آن ها بر روی سیستم بودیم، و تجربه های مفیدی در استفاده از ویژگی های سیستم به دست آوردیم که تاکنون موقعیت استفاده از آن ها را نداشتیم موضوع با اهمیت تر این که گردآوری این رویداد نشان داد که WCT با موفقیت می تواند با تقاضاهای فعالیت گردآوری وب و اندازه افزایش یافته تطبیق دهد مناطقی را که به بررسی و نظارت دقیق تر نیاز داشتند نیز مشخص کند و درس های ارزشمندی آموخته و تجربه ای به دست دهد.

متعاقب انتخابات هیئت محلی سال 2007، رویداد دیگری که بر پایه گردآوری قرار داشت جام جهانی راگی 2007 بود این رویداد در مقیاس بسیار کوچک تر از گردآوری رویداد اول بود، اما شامل وب گاه های خارج از محدوده مقررات واسپاری قانونی زلاندنو می شد که امکان قانونی جمع آوری با استفاده از WCT را برای اولین بار نیاز داشت.

انتخاب گران الکترونیکی کتابخانه به طور روز مره برای انتخاب نمودار گردآوری، و مرور وب گاه ها از WCT استفاده می کنند سپس آن ها را برای آرشیو دیجیتال پیشنهاد می دهند.

انتخاب گران الکترونیکی کار را با ویرایش 1.1 شروع کردند و پس از آشنایی بیشتر با این نرم افزار، از نوعی محدودیت روزافزون آگاه شدند که Specification ویرایش جدیدی را اطلاع می داد. نرم افزارهای کنترل کیفی به روز شده تفاوت عظیمی را به وجود آوردند و بسیاری از سایت ها که برای مرور آن ها با ویرایش های اولیه مشکل بودند، جمع آوری و با ویرایش 1.2 مرور شدند.

هنوز برخی مسائل پر دردسر و آزار دهنده وجود دارد و یا در حال ضبط آن ها در bugtraker روی وب گاه های منبع باز هم ضبط کرده ایم.

از آخرین هفته ها ویرایش جدید نرم افزار (ویرایش 3) قابل بهره برداری شد و انتظار داریم که بتواند در گردآوری گردش کار، نتایج بهتری ایجاد کند.

پس از دو دوره گردآوری کتابخانه، ظرفیت مدل ها برای ذخیره آتی را توسعه داده و نیاز به پهنای باند دارد. ما همچنین دسترسی نرم افزارها و سطح دامنه گردآوری وب گاه های زلاندنورا طراحی می کنیم.

منابع

1. HTTrack Website Copier, (Accessed 2008-05-15).

2. Paynter, G. W. and Mason I. B. (2006) Building a Web Curator Tool for The National Library of New Zealand. New Zealand Library Association (LIANZA) Conference, . October 2006. (Accessed 2008-05-15).

Web Information Management

:Vol. 1

Basics and Worldwide Initiatives

Edited by

Gholam Ali Montazer (Ph.D)

and

Tarbiat Modares University

Farzaneh Shadanpour

National Library and Archives of Islamic Republic of Iran

ص: 318

بسمه تعالی

جَاهِدُوا بِأَمْوَالِكُمْ وَأَنْفُسِكُمْ فِي سَبِيلِ اللَّهِ ذَلِكُمْ خَيْرٌ لَّكُمْ إِنْ كُنْتُمْ تَعْلَمُونَ

با اموال و جان های خود، در راه خدا جهاد نمایید، این برای شما بهتر است اگر بدانید.

(توبه : 41)

چند سالی است که مرکز تحقیقات رایانه ای قائمیه موفق به تولید نرم افزارهای تلفن همراه، کتاب خانه های دیجیتالی و عرضه آن به صورت رایگان شده است. این مرکز کاملاً مردمی بوده و با هدایا و نذورات و موقوفات و تخصیص سهم مبارک امام علیه السلام پشتیبانی می شود.

برای خدمت رسانی بیشتر شما هم می توانید در هر کجا که هستید به جمع افراد خیراندیش مرکز بپیوندید.

آیا می دانید هر پولی لایق خرج شدن در راه اهلبیت علیهم السلام نیست؟

و هر شخصی این توفیق را نخواهد داشت؟

به شما تبریک میگوئیم.

شماره کارت :

6104-3388-0008-7732

شماره حساب بانک ملت :

9586839652

شماره حساب شبا :

IR390120020000009586839652

به نام : (موسسه تحقیقات رایانه ای قائمیه)

مبالغ هدیه خود را واریز نمایید.

آدرس دفتر مرکزی:

اصفهان - خیابان عبدالرزاق - بازارچه حاج محمد جعفر آبا ده ای - کوچه شهید محمد حسن توکلی - پلاک 129/34 - طبقه اول

وب سایت: www.ghbook.ir

ایمیل: Info@ghbook.ir

تلفن دفتر مرکزی: 03134490125

دفتر تهران: 021 - 88318722

بازرگانی و فروش: 09132000109

امور کاربران: 09132000109



مرکز تحقیقات رایانگی

اصفهان

گامی

WWW



برای داشتن کتابخانه های تخصصی
دیگر به سایت این مرکز به نشانی

www.Ghaemiyeh.com

www.Ghaemiyeh.net

www.Ghaemiyeh.org

www.Ghaemiyeh.ir

مراجعه و برای سفارش با ما تماس بگیرید.

۰۹۱۳ ۲۰۰۰ ۱۰۹

