

به نام خدا



پژوهشگاه ارتباطات و فناوری اطلاعات

مطالعه تطبیقی سامانه تحلیل

و جستجوی خبر

پروژه: مدیریت طرح

کد پروژه:

مجری: علیرضا یاری

تهیه کننده: تیم مدیریت طرح

کد گزارش: P-PD-VAS-SBM-S-011-1.02

تاریخ ارائه: ۹۴/۰۵/۱۷

نسخه/وضعیت: نهایی



در راستای تحقق مأموریت پژوهشگاه ارتباطات و فناوری در فراهم سازی سکویی برای ارتقاء دانش، انتقال فناوری و بومی سازی محصولات و خدمات حوزه فاوا و با هدف جلب مشارکت علاقه مندان در توسعه و بهره مندی از دستاوردهای پژوهشگاه ارتباطات و فناوری اطلاعات، آزاد رسانی این دستاوردها در زمره برنامه های اولویت دار پژوهشگاه به شمار می آید. به همین منظور مستند حاضر تحت مجوز بین المللی **CC-BY-SA-NC** نسخه ۴، در دسترس عموم قرار گرفته است. شایان ذکر است تحت این مجوز، ضمن حفظ مالکیت فکری این مستند برای پژوهشگاه ارتباطات و فناوری اطلاعات، بازانتشار و بکارگیری آن صرفاً برای موارد تحقیقاتی و با ذکر نام پژوهشگاه ارتباطات و فناوری اطلاعات بلامانع است.

شناسنامه گزارش

شماره نسخه: ۱,۰۲	عنوان: مطالعه تطبیقی سامانه تحلیل و جستجوی خبر	
تاریخ ارائه گزارش: ۹۴/۰۵/۱۷	نوع گزارش: فنی	کد: P-PD-VAS-SBM-S-011-1.02
نام پروژه: مدیریت طرح	نوع پروژه: پژوهشی - کاربردی	
تاریخ شروع: ۹۳/۱۱/۲۰	تاریخ پایان: ۹۶/۱۱/۲۰ (۳۶ ماه)	
نام گروه: طرح جویسگر		
کد پروژه: ۹۳۳۲۰۱۲	شماره و تاریخ قرارداد: ۹۳/۱۱/۲۰	
مجری: علیرضا یاری	ناظر / ناظرین: کامبیز بدیع، رامین شکری پور و روح‌اله رحمانی	
تهیه کننده / تهیه کنندگان: تیم مدیریت طرح		
نام و نشانی مجری:		
تهران، انتهای خیابان کارگر شمالی، مرکز تحقیقات مخابرات ایران - کد پستی: ۱۴۳۹۹۵۵۴۷۱ - تلفن: ۸۸۰۰۵۵۰۸-۱۰		
نام و نشانی حمایت کننده:		
تهران، خیابان شریعتی، وزارت ارتباطات و فناوری اطلاعات		
ملاحظات: ندارد		
چکیده:		
<p>تحلیل ۵۰ سایت برتر ایران در سایت الکسا حاکی از این است که حدود ۳۰ درصد از این ۵۰ سایت پرترافیک در ایران سایت‌های خبری هستند. این موضوع بیانگر میزان اهمیت این سایت‌های خبری و استفاده زیاد کاربران ایرانی از این سایت‌ها است. لذا یکی از سرویس‌های جذاب و پرکاربرد برای کاربران که می‌تواند یک مزیت رقابتی برای جویسگر بومی در مقایسه با جویسگرهای محبوبی مانند گوگل باشد، خدمات مبتنی بر اخبار می‌باشد. با توجه به وجود سرویس‌دهنده‌گان خبری متعدد خارجی و داخلی و همچنین حجم زیاد اخبار مربوط به موضوعات خبری مختلف لذا برای پاسخگویی به نیازهای اطلاعاتی کاربران در رابطه با اخبار نیازمند طراحی و پیاده‌سازی سامانه‌های تحلیل و جستجوی خبری می‌باشیم. در این مستند هدف مطالعه و بررسی سامانه‌های تحلیل و جستجوی خبری در سطح کشور و دنیا است که در نهایت به دید مناسبی از وضعیت موجود رسید و با توجه به نقاط ضعف و قوت سامانه‌های موجود در کشور بتوان آن‌ها را توسعه و یا بهبود داد.</p>		
کلمات کلیدی: جویسگر بومی		
وضعیت گزارش: نهایی	زبان گزارش: فارسی	
وضعیت دسترسی: عادی	تعداد صفحات: ۸۳	

چکیده

تحلیل ۵۰ سایت برتر ایران در سایت الکسا حاکی از این است که حدود ۳۰ درصد از این ۵۰ سایت پرتراфик در ایران سایت‌های خبری هستند. این موضوع بیانگر میزان اهمیت این سایت‌های خبری و استفاده زیاد کاربران ایرانی از این سایت‌ها است. لذا یکی از سرویس‌های جذاب و پرکاربرد برای کاربران که می‌تواند یک مزیت رقابتی برای جویشر بومی در مقایسه با جویشرهای محبوبی مانند گوگل باشد، خدمات مبتنی بر اخبار می‌باشد. با توجه به وجود سرویس‌دهنده‌گان خبری متعدد خارجی و داخلی و همچنین حجم زیاد اخبار مربوط به موضوعات خبری مختلف لذا برای پاسخگویی به نیازهای اطلاعاتی کاربران در رابطه با اخبار نیازمند طراحی و پیاده‌سازی سامانه‌های تحلیل و جستجوی خبری می‌باشیم. در این مستند هدف مطالعه و بررسی سامانه‌های تحلیل و جستجوی خبری در سطح کشور و دنیا است که در نهایت به دید مناسبی از وضعیت موجود رسید و با توجه به نقاط ضعف و قوت سامانه‌های موجود در کشور بتوان آن‌ها را توسعه و یا بهبود داد.

اطلاعات مرتبط

مستندات مرتبط

شماره مستند	نوع مستند	نام مستند

تغییرات اعمال شده در نسخه‌های پیشین

شماره نسخه	تاریخ	تغییرات اعمال شده
۱،۰	۹۴/۰۵/۱۷	آماده‌سازی گزارش
۱،۰۱	۹۴/۰۶/۱۰	اضافه کردن چکیده، کامل کردن مقدمه، اضافه کردن بخش نتیجه‌گیری، اصلاح شکل‌ها و جداول
۱،۰۲	۹۴/۰۶/۲۲	چکیده و نتیجه‌گیری با توجه به فیدبک ناظرین اصلاح شد

تأییدکنندگان

نام و نام خانوادگی	تاریخ	امضاء	ملاحظات
علیرضا یاری			مجری پروژه
تیم مدیریت طرح			تهیه کننده / تهیه کنندگان
کامبیز بدیع، رامین شکری پور و روح‌اله رحمانی			ناظر پروژه
			مدیر گروه
مانا روزی طلب			مسئول مستندات پژوهشکده
علیرضا یاری			رئیس پژوهشکده / معاون پژوهشی

اسامی اعضای تیم مدیریت طرح بر اساس حروف الفبا

۱. محمد آزادنیا	۱۱. مهدی عمادی
۲. محمد مهدی اثنی اشری	۱۲. مژگان فرهودی
۳. شهره جهانبخش	۱۳. محمد مهدی کیخا
۴. مونا داوودی	۱۴. مریم محمودی
۵. غزاله رحمانی فرزین	۱۵. پویان مسعودی فر
۶. فرزانه رحمانی	۱۶. اکبر مقدر
۷. فرزاد زرگری	۱۷. امین میرزائی
۸. محمد صادق زاهدی	۱۸. محمدرضا میرصراف
۹. علی شریفی	۱۹. حمیدرضا نصیری آسایش
۱۰. معصومه عظیم زاده	۲۰. علیرضا یاری
۱۱. طاهره علوی زرگر	

سرفصل مطالب

سامانه تحلیل و جستجوی خبر

۶

- | | |
|----|--|
| ۷ | ۱-۱ موتور جستجوی خبر بومی |
| ۹ | ۲-۱ موتور جستجوی خبر جهانی |
| ۱۱ | ۳-۱ شاخصهای ارزیابی سامانه تحلیل و جستجوی خبر |
| ۱۱ | ۱-۳-۱ معیارهای ارزیابی خزشگر سامانه تحلیل و جستجوی |
| ۱۱ | ۲-۳-۱ معیارهای ارزیابی بخش پارسر و استخراج اطلاعات |
| ۱۲ | ۳-۳-۱ معیارهای ارزیابی بخش جستجوی خبر |
| ۱۳ | ۴-۳-۱ معیارهای ارزیابی بخش تحلیل خبر |
| ۱۵ | ۱-۳-۵ معیارهای ارزیابی نیازهای غیرکارکردی |

فهرست تصاویر

تصویر ۱: سامانه تحلیل و جستجوی خبر ۷

فهرست جداول

جدول ۱: موتور جستجوی خبر بومی ۷

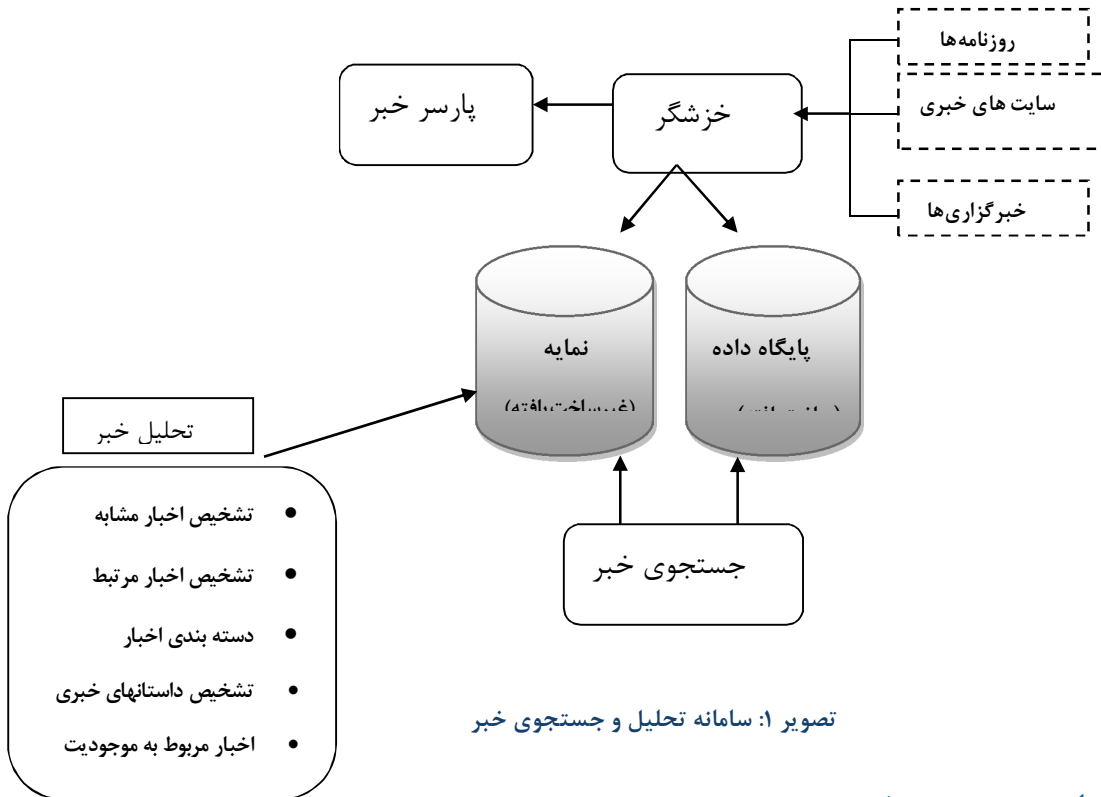
جدول ۲: لیست مرتب شده سایتهای خبری ۹

جدول ۳: فهرست معروفترین موتورهای جستجوی خبر ۱۰

سامانه تحلیل و جستجوی خبر

با توجه به وجود سرویس‌دهنده‌گان خبری متعدد خارجی و داخلی و همچنین حجم زیاد اخبار مربوط به موضوعات خبری مختلف لذا برای پاسخگویی به نیازهای اطلاعاتی کاربران در رابطه با اخبار نیازمند طراحی و پیاده‌سازی سامانه‌های تحلیل و جستجوی خبری می‌باشیم. در این مستند به مطالعه و بررسی سامانه‌های تحلیل و جستجوی خبری در سطح کشور و دنیا پرداخته می‌شود. در ابتدا به بررسی معماری کلی سامانه تحلیل و جستجوی خبر پرداخته و بخش‌های مختلف آن معرفی می‌شود. در ادامه به بررسی سایت‌هایی که در کشور خدمات مبتنی بر اخبار ارائه می‌دهند پرداخته می‌شود. ویژگی‌ها و مشخصه‌های اصلی این سایت‌ها مقایسه و تحلیل شده است و بر این اساس نقاط ضعف و قوت این خدمات دهندگان خبری بومی استخراج شده است. علاوه بر سایت‌های خبری بومی، سامانه‌های خبری جهانی نیز بررسی و مقایسه شده است در انتها نیز معیارهای ارزیابی بخش‌های مختلف سامانه تحلیل و جستجوی خبری بررسی شده است.

در شکل ۱ معماری کلی سامانه تحلیل و جستجوی خبری نشان داده شده است. سامانه تحلیل و جستجوی خبر از مولفه‌های خزشگر، پارسر، نمایه و مولفه‌ی تحلیل و جستجوی خبر تشکیل شده است. در مؤلفه خزشگر خبر، صفحات خبری مختلف از سایت‌های خبری، روزنامه‌ها و خبرگزاری‌های معتبر جمع‌آوری می‌شود. سپس مؤلفه پارسر خبر با استفاده از روش‌ها و الگوریتم‌های مختلف پردازش متن، به استخراج بخش‌های مختلف صفحات خبری می‌پردازد. در این مؤلفه از ابزارهای مختلفی مانند واحدساز، یکسان‌ساز، برچسب‌گذار اجزای کلام، شناسایی کلمات مرکب، ریشه‌یاب و شناسایی جملات فارسی استفاده شده و متون برای تحلیل‌های سطح بالاتر آماده می‌شوند. در بخش تحلیل خبر با استفاده از روش‌های مختلف متن کاوی دسته خبر، اخبار مهم، کلمات کلیدی خبر، خبرهای مشابه، وقایع و داستان‌های خبری استخراج می‌شوند. در مؤلفه نمایش خبر نیز اخبار با واسط کاربری مناسب، به کاربرهای مختلف نمایش داده شده و امکان جستجوی اخبار برای کاربران نیز فراهم شده است.



۱-۱ موتورهای جستجوی خبر بومی

در ابتدا لیستی از مهمترین ویژگی‌هایی که یک سامانه تحلیل و جستجوی خبری باید داشته باشد استخراج گردیده و سپس بر اساس این ویژگی‌ها به مقایسه و بررسی سایت‌های خبری موجود در کشور پرداخته شده است. جدول زیر فهرست برخی از معروف‌ترین وب‌سایت‌های خبری به همراه ویژگی‌های اصلی آنها را نشان می‌دهد. برای هر ویژگی امتیاز از ۰ تا ۵ در نظر گرفته شده و بر اساس امکانات این سایتها تکمیل شده است. برای تکمیل از داده‌های شرکت‌های که در طرح شرکت کرده‌اند نیز استفاده شده است.

جدول ۱: موتور جستجوی خبر بومی

ردیف	ویژگی	خبر فارسی	شهر خبر	خبر پو	قطره	تی نیوز	خبر خون	خبر پارسی جو	خبر یوز
۱	رتبه در الکسا	۴۹	۲۳	۴۶	۶۴	۷۳	۱۲۳	۴۷۹	۶۶۳
							۲		

۰	۳	۳	۲	۰	۰	۰	۴,۵	تشخیص اخبار مشابه(نقل قول ها، اخبار تکراری)	۲
۳	۴	۴	۰	۰	۰	۰	۴	تشخیص اخبار مرتبط(اخبار مرتبط با خبر جاری کاربر)	۳
۱	۵	۲	۴,۵	۰	۳,۵	۴,۵	۲	نمایش لحظه ای اخبار(به صورت زنده)	۴
۲	۴,۵	۳	۴	۰	۴	۵	۴,۵	نمایش پر بیننده ترین اخبار	۵
۰	۰	۲	۴	۰	۴	۰	۰	نمایش پر بحث ترین اخبار	۶
۰	۰	۰	۰	۰	۰	۰	۰	امکان مربوط به کاربران(امکاناتی مانند بحث کردن، دوستی، مشاهده اخبار دوستان و ...)	۷
۰	۰	۲	۰	۰	۰,۵	۰	۱,۵	امکان شخصی سازی اخبار مشاهده شده برای کاربر(بر اساس فیلدهای مختلف نوع خبر، مرجع خبر، موضوع خبر، موجودیت های مختلف مانند افراد، سازمان، مکان ها و ... به صورتی که کاربر اخبار را به صورت فیلتر شده بر اساس پروفایل خود ببیند)	۸
۲,۵	۴	۳,۵	۴	۲,۵	۲,۵	۳	۵	نمایش خبر اصلی (امکانات و نحوه ی نمایش خبر موقع کلیک کاربر بر روی خبر)	۹
۳,۵	۳,۵	۲	۳,۵	۲	۱	۱,۵	۳	اخبار چند رسانه ای (تصویر، ویدئو، ...)	۱۰
۴	۵	۳	۳	۱	۲,۵	۲,۵	۴	واسط کاربری سایت(صفحه ی اول سایت، طراحی سایت و ...)	۱۱
۳,۵	۳,۵	۲	۳,۵	۲	۰	۰	۳,۵	امکانات جستجو (دقت جستجو، جستجوی پیشرفته)	۱۲
۴	۴	۴	۴	۴	۴	۴	۴	دسته بندی اخبار(بر اساس فیلدهای مختلف مانند موضوع، استان، کشور و ...)	۱۳
۲	۳	۱,۵	۰	۰	۱	۰	۲,۵	تشخیص داستانهای خبری (news Stories) و امکانات مربوط به نمایش آن مانند نمایش پر بیننده ترین، پر بحث ترین داستان ها، دسته بندی موضوعی داستان های خبری و ...	۱۴
۲	۰	۰	۱	۰	۰	۰	۰	اخبار مربوط به موجودیت ها(مانند اشخاص، سازمان و ...)	۱۵
۳	۲	۴	۴	۵	۴,۵	۳	۵	تنوع سایت های خبری پوشش داده شده	۱۶

۲	۲	۱,۵	۱,۵	۰	۰	۰	۲	امکانات تحلیل اخبار شامل مواردی مانند (شناسایی موجودیت ها در اخبار مانند مکان ها ، افراد، سازمان ها و تحلیل های مرتبط با آنها، تشخیص شبکه انتشار اخبار، کاوش نظرات کاربران در مورد اخبار مربوط به موجودیت های مختلف مانند سازمان یا اشخاص) لطفا در مورد سیستم خبر پیشنهادی خود سایر تحلیل ها در ستون مربوطه آورده شود	۱ ۷
۲۲,۵	۴۳,۵	۳۷,۵	۳۹	۱۶,۵	۲۷,۵	۲۳,۵	۴۵,۵		

بر اساس امتیازهای بالا لیست مرتب شده سایتهای خبری به صورت زیر است:

جدول ۲: لیست مرتب شده سایتهای خبری

رتبه	نام موتور جستجو	امتیاز	رتبه الکسا
۱	خبر فارسی	۴۵,۵	۴۹
۲	خبر پارسی جو	۴۳,۵	۴۷۹
۳	تی نیوز	۳۹	۷۳
۴	خبر خون	۳۷,۵	۱۲۳۲
۵	خبر یوز	۳۲,۵	۶۶۳

با بررسی و تحلیل مطالعه صورت گرفته مشخص گردید که ضعفهای اصلی این سایتها بیشتر مربوط به بخش تحلیل اخبار و شخصی سازی اخبار و تشخیص وقایع و داستان های خبری است. در واقع سایت های که در کشور خدمات مربوط به خبر ارائه می دهند بیشتر در قالب تجمیع کننده اخبار هستند و خدماتی مانند تحلیل اخبار و یا شخصی سازی اخبار را ارائه نمی دهند. با توجه به اهمیت سرویس های خبری موتورهای جستجو اصلی موجود در کشور نیز سرویس خبر ارائه می دهند. هر یک از این موتورهای جستجو ویژگی های منحصر به فرد دارند به عنوان مثال سرویس خبری یوز امکان مشاهده پدیده های خبری را برای کاربر فراهم نموده است و یا یک بخش مربوط به کارکاتورهای خبری دارد که ممکن است برای کاربر جذاب باشد همچنین خبرپارسی جو نقاط قوتی مانند واسط کاربری مناسب، امکان مشاهده زنده اخبار، پوشش خبری (جمع آوری و دسته بندی اخبار مشابه) و نمایش کلیه اخبار متنی و چند رسانه ای در یک صفحه، دارد. در کل خبرپارسی جو بر اساس معیارهای مشخص شده در جدول ۱ عملکرد بهتری داشته است.

جدول زیر فهرست برخی از معروفترین موتورهای جستجوی خبر به همراه ویژگی‌های اصلی آنها را نشان

می دهد

جدول ۳: فهرست معروفترین موتورهای جستجوی خبر

ویژگی‌ها	رتبه محلی	رتبه جهانی	
<p>نسخه ی اولیه در سال ۲۰۰۲ و نسخه رسمی در سال ۲۰۰۶ شروع شده است ۴۵۰۰ منبع خبری را پوشش داده است نسخه ی محلی برای ۶۰ منطقه و ۲۸ زبان توسعه داده شده است(فارسی ندارد) ۲۵۰۰۰ منتشر کننده خبر برای زبان انگلیسی و ۴۵۰۰ برای سایر زبانها امکان جستجوی خبر و مرتب سازی نتایج بر اساس زمان انتشار خبرها دسته بندی موضوعی خبرها امکان تعریف Alerts برای کلمات کلیدی کاربر و ارسال اخبار مرتبط با آن کلمات بای ایمیل کاربر امکان شخصی سازی اخبار توسط کاربر وجود نسخه ی موبایلی اخبارگوگل امکان دریافت اخبار از طریق فید یا RSS امکان جستجو در آرشیو اخبار مربوط به ۲۰۰ سال گذشته امکان جستجوی پیشرفته بر اساس فیلدهای مختلف مانند خبرگزاری، زمان انتشار خبر، عنوان خبر، بدنه ی خبر ، مکان و ... امکان مشاهده داستان های خبری پربیننده (دسته‌ای از رویدادهای خبری با یکدیگر یک داستان خبری را می‌سازند. مثلاً «بازی‌های آسیایی» یک داستان خبری است. این داستان هر چهارسال یکبار تکرار می‌شود و در زمانی که در حال انجام است بیش‌ترین اخبار را به خود اختصاص می‌دهد. اخبار مربوط به این داستان در طول زمان ذخیره‌شده و کاربر می‌تواند تمامی اخبار مربوط به یک داستان را به صورت یکجا مشاهده کند.) نمایش اخبار مرتبط با خبر جستجو شده حذف اخبار تکراری نمایش اخبار بر اساس اخبار مشاهده شده کاربر نمایش اخبار پر بازدید خود خبر اصلی در سایت اصلی منتشر کننده خبر نمایش داده می شود در نتایج گوگل فقط خلاصه خبر نمایش داده می شود علاوه بر آن برای هر عنوان خبری اخبار مربوط به سایتهای مختلف منتشر کننده دسته بندی شده است بر اساس موارد مختلفی مانند : نظر شخصی(Opinion)، تحلیلی(In Depth)، تعداد ارجاع ها (Highly Cited) نمایش داده شده است، کلیه موارد دیگر برای عنوان خبری مانند عکس ها و ویدئو دسته بندی شده و به صورت مناسبی نمایش داده شده است. نمایش خط زمانی خبر به صورت نمودار</p>	۱	۱	Google News
<p>رتبه دوم در سایت الکسا در رتبه بندی مربوط به سایت های خبری دسته بندی اخبار و امکان نمایش اخبار هر دست متناسب با آن به عنوان مثال دسته ی تکنولوژی به صورت عکس نمایش داده شده است یک عکس مرتبط با خبر که بخشی از خبر نیز در آن نمایان است. و در مان بخش تکنولوژی دسته بندی های دیگری مانند خبرهای مربوط به review رکت اپل و ... یکی از استراتژی های اصلی یاهو تمرکز بر تولید محتوای اصلی(original content) امکان جستجوی خبر</p>	۵	۵	Yahoo! News

رتبه اول در سایت الکسا در رتبه بندی مربوط به سایت های خبری شبکه اجتماعی خبر است که در آن کاربران ثبت نام کرده قادرند اخبار را در قالب لینک یا متن ارسال کنند و آن را با دیگران به اشتراک بگذارند. افراد می توانند گروه تشکیل دهند. یکی از بخش های معروف ردیت بخش «من هستم» است که در آن کاربران اعلام می کنند که هر چه می خواهید از من بپرسید تاکنون بسیاری از افراد مشهور نظیر، باراک اوباما، جیمی کیمل، ران پال، استیون کلبر، ددماو، زاک براف و نیل استراس در این بخش شرکت کرده اند کاربران می توانند به پست یا گزارش های نوشته شده رأی دهند و جایگاه آن خبر را در صفحات ردیت و صفحه اصلی مشخص کنند. در واقع کاربران هستند که در ردیت اهمیت خبر را تعیین می کنند	۱۰	۳۱	Reddit
تعداد کارمندان آن ۷۱ نفر است از سال ۲۰۰۵ راه اندازی شده است تعداد بازدیدکننده ردیت در ماه ژوئیه ۲۰۱۵ حدود ۱۹۵،۲۰۹،۱۶۹ از ۲۱۵ کشور مختلف که تعداد ۸،۰۶۴،۶۷۸،۱۱۲ صفحه را مشاهده کرد اند. امکان جستجوی پیشرفته بر اساس نویسنده خبر، سایت و ... ثبت نام نیازی به ایمیل برای تکمیل ندارد امکان شخصی سازی برای کاربر امکان گذاشتن یک پست خبری بدون رفرنس خارجی تحت عنوان self posts با پایتون توسعه داده شده است			

۳-۱ شاخص های ارزیابی سامانه تحلیل و جستجوی خبر

هر یک از بخش های سامانه تحلیل و جستجوی خبر را می توان به صورت مجزا بر اساس معیارهای مختلف ارزیابی نمود. در ادامه برخی از معیارها و آزمون های ارزیابی بخش های مختلف توضیح داده خواهد شد.

۱-۳-۱ معیارهای ارزیابی خزشگر سامانه تحلیل و جستجوی خبری

- ۱- سرعت: تعداد و حجم خبرهای جمع آوری شده در واحد زمان
- ۲- پوشش خبر و تازه بودن: در این معیار تعداد خبرگزاری ها، روزنامه ها و سایتهای خبری پوشش داده شده مورد بررسی قرار می گیرد و یک مدت زمان ماکزیمم برای فاصله انتشار آگهی و جمع آوری آن توسط خزشگر تخمین زده می شود. بدیهی است هر چقدر زمان تاخیر کمتر باشد کارایی خزشگر بیشتر است.
- ۳- تحمل پذیری خطا: در این بخش از ارزیابی، رفتار خزشگر در مواجهه با انواع خطاهای شبکه، دامها حلقه بررسی می شود.

۲-۳-۱ معیارهای ارزیابی بخش پارسر و استخراج اطلاعات

مهمترین معیارها برای سنجش کارایی این ماژول معیارهای دقت (precision)، بازآوری (Recall) و معیار F است. که به صورت زیر تعریف می‌شوند:

$$\text{Precision} = \frac{NNN \ NNNN \ NN \ NNNNNNNNN \ NNNNNNNNNNN \ NNNNN}{NNN \ NNN \ NN \ NNNNNNNNN \ NNNNNN \ NNNNN}$$

$$NNNNNN = \frac{NNN \ NNNN \ NN \ NNNNNNNNN \ NNNNNNNNNNN \ NNNNN}{NNN \ NNN \ NN \ NNNNNN \ NNNNN}$$

معیار F که میانگین هارمونی دقت و بازآوری است نیز به صورت زیر تعریف می‌شود:

$$N = \frac{2PR}{P+R}$$

۱-۳-۳ معیارهای ارزیابی بخش جستجوی خبر

معیار دقت در مرتبه n یا P@N: این معیار کسری از خبرهای مرتبط به پرس‌وجو را که در n رتبه اول ارزیابی شده‌اند نشان می‌دهد. مقادیر مرسوم برای n عبارتند از: ۵، ۱۰، ۲۰.

معیار متوسط میانگین دقت (MAP): میانگین دقت (average precision) در مرتبه n را برای مقادیر مختلف n برای یک پرس‌وجو داده شده نشان می‌دهد. متوسط میانگین دقت را برای یک مجموعه پرس‌وجو (حداقل ۵۰ پرس‌وجو) نان می‌دهد.

$$\text{MAP}(Q) = \frac{1}{|Q|} \sum_{N \in Q} \frac{1}{m_{NN}} \text{Precision}(R_N)$$

معیار NDCG: در این معیار، علاوه بر اینکه مرتبط بودن هر خبر مشخص می‌شود، رتبه آن نیز در عمل مقایسه، تاثیرگذار است؛ از آنجایی که کاربران علاقمند به رؤیت صفحات مرتبط در بالای لیست هستند، در این روش نتایجی که ارتباط بیشتری با پرس‌وجوی کاربر دارند در صورتی که رتبه بالاتری داشته باشند از اهمیت بیشتری در محاسبه برخوردارند. در محاسبه NDCG ابتدا باید مقدار CG را محاسبه کنیم برای این منظور مجموع میزان ارتباط هر نتیجه ارائه شده در لیست نتایج با پرس‌وجوی کاربر را به دست می‌آوریم.

$$\text{CG}_N = N \sum_{N \in Q} \text{rel}_N$$

که در اینجا rel_i میزان ارتباط نتیجه i ام با پرس‌وجو می‌باشد. عددی که به کمک این تابع به دست می‌آید رتبه‌بندی سیستم را مد نظر قرار نمی‌دهد و تنها معیاری برای مرتبط بودن هر سند است، بنابراین اگر نتیجه‌ای که ارتباط کمتری با پرس‌وجو دارد در رتبه بالاتری نسبت به نتیجه مرتبط‌تر قرار بگیرد در مقدار CG تفاوتی ایجاد نمی‌شود. با توجه به این مشکل از DCG برای لحاظ نمودن رتبه‌بندی استفاده می‌کنیم.

$$\text{DCG}_N = \text{rel}_1 + N \sum_{N \in Q} \frac{\text{rel}_N}{\log_2 i}$$

به این ترتیب اگر نتایجی که ارتباط بیشتری دارند در رتبه پایین تر قرار بگیرند سودمندی و کارایی سیستم کاهش می‌یابد. با توجه به اینکه اندازه نتایج جستجو برای پرس‌وجوهای مختلف متفاوت است بنابراین برای مقایسه کارایی موتورهای جستجو باید DCG را نرمال کنیم.

$$nDCG_N = \frac{DCG_N}{IDCG_N}$$

در اینجا IDCG_p همان DCG ایده‌آل در نقطه p است.

۱-۳-۴ معیارهای ارزیابی بخش تحلیل خبر

هر یک از بخشهای تحلیل خبر ارزیابی خاص خودش را نیاز دارد. در ادامه معیارهای مربوط به برخی از این بخش‌ها آورده می‌شود:

بخش خوشه بندی اخبار: روش‌های ارزیابی خوشه‌های حاصل از خوشه‌بندی را به صورت سه دسته تقسیم می‌کنند که عبارتند از:

- معیارهای خارجی (External Criteria)
- معیارهای درونی (Internal Criteria)
- معیارهای نسبی (Relative Criteria)

هم معیارهای خارجی و هم معیارهای درونی بر مبنای روش‌های آماری عمل می‌کنند و پیچیدگی محاسباتی بالایی را نیز دارا هستند. معیارهای خارجی عمل ارزیابی خوشه‌ها را با استفاده از بینش خاص کاربر انجام می‌دهند. معیارهای درونی عمل ارزیابی خوشه‌ها را با استفاده از مقادیری که از خوشه‌ها و نمای آنها محاسبه می‌شود، انجام می‌دهند.

پایه معیارهای نسبی، مقایسه بین شمای خوشه‌بندی (الگوریتم به علاوه پارامترهای آن) مختلف است. یک و یا چندین روش مختلف خوشه‌بندی چندین بار با پارامترهای مختلف روی یک مجموعه داده اجرا می‌شوند و بهترین شمای خوشه‌بندی از بین تمام شمای انتخاب می‌شود. در این روش مبنای مقایسه، شاخص‌های اعتبارسنجی (Validity-Index) هستند.

شاخص‌های اعتبارسنجی برای سنجش میزان صحت (Goodness) نتایج خوشه‌بندی به منظور مقایسه بین روش‌های خوشه‌بندی مختلف یا مقایسه نتایج حاصل از یک روش با پارامترهای مختلف مورد استفاده قرار می‌گیرند مانند:

- Davies-Bouldin
- Calinski-Harabasz
- Dunn index
- R-square index
- Hubert-Levin(C-index)

- Krzanowski-Lai
- Hartigan index
- Root-mean-square standard deviation (RMSSTD) index
- Semi-partial R-squared (SPR) index
- Distance between two clusters (CD) index
- weighted inter-intra index
- Homogeneity index
- Separation index

بخش خلاصه‌سازی اخبار:

Precision و recall به عنوان معیارهای متداول در خلاصه‌سازی مطرح شده است. recall برابر است با تقسیم تعداد جملاتی که توسط سیستم درست تشخیص داده شده است بر تعداد جملاتی که توسط انسان درست تشخیص داده شده اند. Precision برابر است با تقسیم تعداد جملاتی که توسط سیستم درست تشخیص داده شده‌اند بر تعداد کل جملاتی که توسط سیستم برای خلاصه ایجاد شده‌اند.

$$Precision = \frac{SentenceCount(Summary_{SystemExtracted} \cap Summary_{Ideal})}{SentenceCount(Summary_{SystemExtracted})}$$

$$Recall = \frac{SentenceCount(Summary_{SystemExtracted} \cap Summary_{Ideal})}{SentenceCount(Summary_{Ideal})}$$

مشخص است که هر چقدر طول خلاصه افزایش یابد، تفاوت بین متن خلاصه شده توسط انسان و متن خلاصه شده اتوماتیک افزایش می‌یابد. در حالت کلی ارزیابی به دو گروه اصلی و فرعی تقسیم می‌شود. در ارزیابی اصلی، کیفیت خلاصه سازی از طریق مقایسه آن با خلاصه های ایده آل دستی انجام می‌شود. در ارزیابی فرعی مشخص می‌شود خلاصه ایجاد شده چقدر روی انجام کارهای دیگر تاثیر گذار بوده است. از جمله این کارها میتوان به پرسش و پاسخ، درک مطلب و مرتبط بودن سند به موضوع خاص نام برد. در صورتی که خلاصه ایجاد شده از نوع چکیده باشد ارزیابی از طریق مقایسه محتوایی با خلاصه دست ساز انسان می‌باشد (ارزیابی محتوایی). این نوع ارزیابی نیاز به تحلیل زبان شناختی دارد. مثلا فاصله بین بردار تکرار کلمات خلاصه ایجاد شده توسط انسان و سیستم محاسبه می‌شود. نوع دیگر ارزیابی، ارزیابی موضوعی می‌باشد. در این حالت از افراد خواسته میشود که خلاصه سیستم را از دو دیدگاه ارزیابی و امتیاز دهی نمایند. اول اینکه خلاصه سیستم به چه میزان از محتویات مهم مقاله اصلی را می‌پوشاند و دوم آنکه خلاصه سیستم تا چه اندازه خوانا می‌باشد.

شناسایی موجودیت ها:

معیارهای دقت (precision)، بازآوری (Recall) و معیار F مهمترین معیارهای به صورت کلی برای سیستم های استخراج اطلاعات هستند. که در این بخش قابل استفاده هستند.

شناسایی وقایع و داستان های خبری:

موارد مختلفی در پژوهش های قبلی انجام شده است مانند موارد زیر.

۲,۲ Event Tracking

۲,۳ Event Summarization (What happened)

در اکثر این پژوهش‌ها از معیارهای دقت (precision)، بازآوری (Recall) استفاده شده است که در هر بخش متناسب با آن تعریف می‌شود.

۱-۳-۵ معیارهای ارزیابی نیازهای غیرکارکردی

- دسترس‌پذیری^۱
 - میزان پاسخگویی سیستم در قبال پرس‌وجوهای ارسال شده
 - نداشتن شکست در سیستم
- کارایی^۲
 - زمان پاسخگویی به پرس‌وجوها
 - توانایی در پاسخگویی به کاربر همزمان
- آزمون گرافیک واسط کاربری^۳
- قابلیت همکاری^۴
 - داشتن وب‌سرویس مناسب

۱-۴ نتیجه‌گیری

با توجه اهمیت سرویس تحلیل و جستجوی خبر تمامی موتورهای جستجو در کشورهای مختلف سرویس خبر ارائه داده‌اند. تحلیل سایت‌های پرتراфик کشور نیز حاکی از این است که سایت‌های خبری جزو سایت‌های پربازدید در کشور است. لذا با توجه به تعداد زیاد خبرگزاری‌ها و سایت‌های خبری و همچنین انتشار حجم عظیمی از انواع اخبار از دسته‌های موضوعی مختلف برای پاسخگویی به نیازهای اطلاعاتی کاربران نیازمند طراحی راه‌کارهای مناسب برای تسهیل کاربران در پیدا کردن اطلاعات مورد نظر خود خواهیم بود. بررسی سایت‌های خبری موجود در کشور نشان می‌دهد که اکثر این سایت‌ها به صورت تجمیع‌کننده اخبار هستند و خدماتی مانند تحلیل و جستجوی اخبار و یا امکان شخصی‌سازی اخبار برای کاربران فراهم نشده است. لذا نیازمند توسعه‌ی سامانه تحلیل و جستجوی خبر خواهیم بود البته هدف از این سامانه توسعه‌ی خدماتی است که با پوشش سایت‌های خبری و خبرگزاری‌ها و

^۱ Availability

^۲ Performance

^۳ Gui Testing

^۴ Interoperability

روزنامه‌های مختلف امکان انواع جستجوی خبر را برای کاربر فراهم کرده و با تحلیل خودکار اخبار کاربر را در رفع نیازهای اطلاعاتی خود کمک نماید. در این سامانه هدف فقط ایجاد یک تجیمع کننده اخبار نیست که اخبار را از تمام سایتهای خبری استخراج کرده و به کاربر نمایش دهد بلکه سامانه‌ی طراحی شده باید کاربر را در یافتن اخبار مورد نظر خود کمک کرده و با امکاناتی که ارائه می نماید رفع نیازهای اطلاعاتی خبری کاربر را تسهیل نماید. به عنوان مثال با تحلیل خودکار محتوای اخبار از خبرگزارهای مختلف اخبار تکراری را حذف کرده همچنین خبرهای مرتبط با پرس و جوی کاربر را نمایش دهد و یا امکان تحلیل نظرات کاربران در مورد اخبار را داشته باشد.