

پژوهشکده فناوری اطلاعات

گروه سکویهای فناوری اطلاعات و فضای مجازی

گزارش فنی

گزارش آزمون پروژه وردنت فارسی در حوزه فاوا

مستخرج از پروژه: توسعه وردنت عمومی زبان فارسی

کد پروژه: ۸۹۳۲۴۱۶

مجری:

محرم منصوری زاده

تهیه کننده/ تهیه کنندگان:

محرم منصوری زاده، محمد نصیری،

محمد دادرس

کد گزارش:

ITF.ITP.TCH.۸۹۳۲۴۱۶.۴۰.۷.۰۲


تاریخ ارائه:

۹۱/۱۱/۱۱

نسخه/ وضعیت


۲،۰ / نهایی

	عنوان گزارش: گزارش آزمون پروژه وردنت فارسی در حوزه فاوا		پژوهشکده
	وضعیت گزارش: نهایی	کد گزارش: ITF.ITP.TCH.۸۹۳۲۴۱۶.۴۰.۷۰۲	فناوری اطلاعات
شناسنامه گزارش			
عنوان گزارش: گزارش آزمون پروژه وردنت فارسی در حوزه فاوا		شماره نسخه: ۲,۰	
کد: ITF.ITP.TCH.۸۹۳۲۴۱۶.۴۰.۷۰۲		نوع گزارش: فنی	تاریخ ارائه گزارش: ۹۱/۱۱/۱۱
نام پروژه: توسعه وردنت زبان فارسی در حوزه فاوا		نوع پروژه: راهبردی-توسعه ای - بنیادی	
تاریخ شروع: ۹۰/۰۴/۲۸		تاریخ پایان: ۹۱/۱۱/۱۱	
نام گروه: سکویهای فناوری اطلاعات و فضای مجازی			
کد پروژه: ۸۹۳۲۴۱۶		شماره و تاریخ قرارداد: ۹۰/۰۴/۲۸-۵۰/۶۶۱۴/ت	
مجری: محرم منصوری زاده		ناظر / ناظرین: مریم محمودی، مژگان فرهودی، بهروز مینایی بیدگلی	
تهیه کننده / تهیه کنندگان: محرم منصوری زاده، محمد نصیری، محمد دادرس			
نام و نشانی مجری: همدان، دانشگاه بوعلی سینا، دانشکده مهندسی، گروه مهندسی کامپیوتر، کد پستی ۶۵۱۴۸۳۳۶۹۵ تلفن: ۸-۰۸۱۱-۸۲۹۲۵۰۵ فکس: ۰۸۱۱-۸۲۹۲۶۳۱			
نام و نشانی حمایت کننده: تهران، انتهای خیابان کارگر شمالی، پژوهشگاه فضای مجازی، کد پستی ۱۴۳۹۹۵۵۴۷۱ تلفن: ۸۴۹۷۷۷۷			
ملاحظات: ندارد			
چکیده: آزمون پروژه وردنت از جنبه‌های کارکردی، غیر کارکردی و محتوایی انجام می‌گیرد. در این سند نحوه انجام آزمون‌های متعدد را بیان نموده و نتایج حاصل از آن‌ها را گزارش می‌کنیم.			
کلمات کلیدی: نتایج آزمون، آزمون نرم افزار، آزمون‌های کارکردی، آزمون‌های غیر کارکردی، آزمون‌های محتوایی			
وضعیت گزارش: نهایی		زبان گزارش: فارسی	
وضعیت دسترسی: عادی		تعداد صفحات: ۲۵	

	عنوان گزارش: گزارش آزمون پروژه وردنت فارسی در حوزه فاوا		پژوهشکده
	وضعیت گزارش: نهایی	کد گزارش: ITF.ITP.TCH.۱۳۹۹۱۳۶۰.V.۰۲	فناوری اطلاعات

چکیده


آزمون پروژه وردنت از جنبه‌های کارکردی، غیر کارکردی و محتوایی انجام می‌گیرد. در این سند نحوه انجام آزمون‌های متعدد را بیان نموده و نتایج حاصل از آنها را گزارش می‌کنیم.

	عنوان گزارش: گزارش آزمون پروژه وردنت فارسی در حوزه فاوا		پژوهشکده
	وضعیت گزارش: نهایی	کد گزارش: ITF.ITP.TCH.AMIRKAZEM.V.02	فناوری اطلاعات

لیست مستندات مرتبط		
شماره مستند	نوع مستند	نام مستند
۴۰	گزارش فنی	گزارش آزمون وردنت فارسی در حوزه فاوا


لیست تغییرات اعمال شده در نسخه‌های قبلی گزارش		
شماره نسخه	تاریخ	تغییرات اعمال شده
۱,۰	۱۳۹۱/۰۲/۲۶	تهیه نسخه اول
۲,۰	۱۳۹۱/۱۱/۱۱	به روز رسانی آمار بر اساس آخرین ویرایش دادگان

تایید کنندگان				
ملاحظات	امضاء	تاریخ	نام و نام خانوادگی	
			محرم منصوری زاده	مجری پروژه
			محرم منصوری زاده، محمد نصیری، محمد دادرس	تهیه کننده / تهیه کنندگان
			مریم محمودی، مژگان فرهودی، بهروز مینایی بیدگلی	ناظر پروژه
			علیرضا یاری	مدیر گروه
			زهره ساعی	مسئول مستندات پژوهشکده
			علیرضا یاری	رئیس پژوهشکده / معاون پژوهشی

	عنوان گزارش: گزارش آزمون پروژه وردنت فارسی در حوزه فاوا		پژوهشکده
	وضعیت گزارش: نهایی	کد گزارش: ITF.ITP.TCH.۱۳۹۹۱۳۲۰.V.۲	فناوری اطلاعات


تقدیر و تشکر

بدین وسیله از ناظرین و مشاورین محترم پروژه قدردانی می‌شود.


	عنوان گزارش: گزارش آزمون پروژه وردنت فارسی در حوزه فاوا		پژوهشکده
	وضعیت گزارش: نهایی	کد گزارش: V.01.ITF.ITP.TCH.1397.01.01.01	فناوری اطلاعات

فهرست مطالب

شماره صفحه	عنوان
۱	۱- پیشگفتار
۱	۱-۱- گستره آزمون
۱	۱-۲- اشخاص و وظایف آن‌ها
۳	۲- معرفی آزمون‌ها
۳	۱-۲- آزمون نرم افزارها
۳	۲-۲- آزمون محتوا
۴	۱-۲-۲- عامل، ملاک، نشانگر و استاندارد
۵	۲-۲-۲- استانداردهای مطلوبیت
۵	۲-۳- روال‌های ارزیابی دقت، صحت و کیفیت محتوا
۵	۱-۳-۲- آزمون‌های واژگانی
۵	۱-۳-۲-۱- آزمون املا
۶	۲-۳-۲- آزمون آوا
۶	۳-۳-۲- آزمون عضویت در ترادف
۶	۴-۳-۲- ارزیابی دستی واژگان
۶	۲-۳-۲- آزمون ترادف
۶	۳-۳-۲- ارزیابی صحت و دقت رابطه‌های معنایی
۶	۴-۳-۲- آزمون بازنمایی
۸	۳- نتایج آزمون‌ها
۸	۱-۳- آمار وردنت فاوا
۸	۲-۳- آزمون نرم افزارها
۹	۳-۳- آزمون صحت و اعتبار دامنه
۹	۴-۳- آزمون صحت و اعتبار منابع واژگانی
۹	۵-۳- آزمون‌های صحت و اعتبار واژه
۹	۱-۵-۳- تعلق به حوزه فاوا
۱۰	۲-۵-۳- املا
۱۰	۳-۵-۳- آوا
۱۰	۴-۵-۳- افزونگی

	عنوان گزارش: گزارش آزمون پروژه وردنت فارسی در حوزه فاوا		پژوهشکده
	وضعیت گزارش: نهایی	کد گزارش: V.07.ITF.ITP.TCH.1397.01.01	فناوری اطلاعات


۱۰.....	آزمون صحت و اعتبار مترادفها.....	۳-۶-
۱۰.....	آزمون عضویت واژه در مترادف.....	۳-۶-۱-
۱۱.....	آزمون شباهت معنایی واژگان.....	۳-۶-۲-
۱۱.....	آزمون تعریف مترادف.....	۳-۶-۳-
۱۱.....	مقوله دستوری.....	۳-۶-۴-
۱۲.....	دامنه اصلی و فرعی.....	۳-۶-۵-
۱۲.....	صحت و اعتبار رابطه های معنایی.....	۳-۷-
۱۳.....	نگاشت وردنت انگلیس به فارسی.....	۳-۸-
۱۴.....	آزمون کیفی محتوا.....	۳-۹-
۱۴.....	آزمون بازنمایی.....	۳-۹-۱-
۱۴.....	آزمون افزونگی مترادف.....	۳-۹-۲-
۱۵.....	آزمون دور و پل در روابط معنایی.....	۳-۹-۳-
۱۶.....	ارزیابی عوامل و ملاکها.....	۳-۱۰-
۱۷.....	خلاصه نتایج آزمونها.....	۳-۱۱-
۱۹.....	جمع بندی.....	۴-

	عنوان گزارش: گزارش آزمون پروژه وردنت فارسی در حوزه فاوا		پژوهشکده فناوری اطلاعات
	وضعیت گزارش: نهایی	کد گزارش: V.۰۲:ITF.ITP.TCH.۱۳۹۳۱۳۰۰۰۰	

فهرست اختصارات

ICT

Information and Communication Technology

	عنوان گزارش: گزارش آزمون پروژه وردنت فارسی در حوزه فاوا		پژوهشکده
	وضعیت گزارش: نهایی	کد گزارش: ITF.ITP.TCH.AMIRKAZEM.V.01	فناوری اطلاعات

هر کدام از اجزای اصلی تشکیل دهنده محتوا را عامل می‌گوییم. در وردنت فاوا واژه (انگلیسی و فارسی)، ترادف (انگلیسی و فارسی)، رابطه معنایی بین ترادفها و رابطه ساختاری بین واژگان عوامل اصلی موجود در وردنت می‌باشند. ویژگی‌های اصلی عامل‌ها را ملاک می‌گویند. مثلاً املا و میزان کاربرد از ملاک‌های عامل واژه به حساب می‌آیند. جزئی‌ترین ویژگی هر ملاک نشانگر نامیده می‌شود. نشانگرها کمیت‌های قابل اندازه‌گیری هستند که برای گردآوری آماره‌هایی درباره عامل‌ها استفاده می‌شوند. سطح مطلوب و مورد پذیرش نشانگر را استاندارد می‌گویند. عوامل اصلی محتوای وردنت فاوا منابع واژگانی، واژگان، ترادفها و انواع روابط معنایی هستند. پرکاربرد بودن، اعتبار و درستی نیز از جمله مهم‌ترین ملاک‌های این عوامل به شمار می‌روند.

۲-۲-۲ - استانداردهای مطلوبیت

مقادیر برخی از نشانگرها را نمی‌توان به صورت قطعی تعریف کرد. مثلاً در بخش منابع واژگانی استفاده از منابع معتبر یکی از نشانگرهای عامل منابع واژگانی است. اما معتبر بودن یک منبع مشخصه‌ای کیفی است و ممکن است منبعی از دیدگاه یک کارشناس معتبر و از دیدگاه کارشناس دیگر نامعتبر تلقی شود. در این پروژه ما بر اساس تشخیص کارشناسان و مشاوران پروژه مقادیر این نشانگرها را تعریف می‌کنیم. به عنوان نمونه ملاک معتبر بودن یا نبودن منبع، تشخیص مشاوران فناوری اطلاعات ماست.


۲-۳ - روال‌های ارزیابی دقت، صحت و کیفیت محتوا

به صورت عمده روال‌های ارزیابی محتوا به سه دسته خودکار، نیمه خودکار و دستی تقسیم می‌شوند. در روش خودکار آزمون تعریف شده به صورت اتوماتیک انجام گرفته و پاسخ واژگان به آزمون گزارش می‌شوند. در روش نیمه خودکار بخشی از واژگان یا ترادفها انتخاب شده و آزمون‌های خاصی روی این بخش صورت می‌گیرد. در روش دستی بخشی افراد خبره بخشی از واژگان یا ترادفها را انتخاب کرده و دقت و صحت آن‌ها را ارزیابی و اعلام نظر می‌نمایند.

۲-۳-۱ - آزمون‌های واژگانی

اطلاعاتی که درباره واژگان ضبط شده صورت املائی واژه و تلفظ آن است. طی دو آزمون املا و آزمون تلفظ دقت و صحت واژه این دو قلم اطلاعات ارزیابی می‌شوند. یکی از مشکلات ثبت واژگان تفاوت‌های رسم‌الخطی یا همان اختلاف صفحات کد مختلف فارسی است. مثلاً در برخی از کامپیوترها کد استفاده شده برای حرف «ی» ۱۶۱۰ است در حالی که در کامپیوتر دیگری این کد ۱۶۰۹ یا ۱۷۴۰ می‌باشد. گرچه شکل نگارش این حرف در هر سه حالت مشابه هم‌دیگر است، هنگام مقایسه حرف به حرف کلمات این کدها متفاوت در نظر گرفته می‌شوند. در پایگاه داده وردنت از کد واحدی برای نمایش هر حرف استفاده می‌شود و هنگام ورود اطلاعات کدها یکسان سازی می‌شوند.

۲-۳-۱-۱ - آزمون املا

	عنوان گزارش: گزارش آزمون پروژه وردنت فارسی در حوزه فاوا		پژوهشکده
	وضعیت گزارش: نهایی	کد گزارش: ITF.ITP.TCH.AMIRKAZEM.V.01	فناوری اطلاعات

آزمون املا به صورت خودکار واژگانی را استخراج می‌کند که در جدول واژگان مرجع وجود ندارند. از آنجایی که همه واژگان ثبت شده در پایگاه داده وردنت طی فرایند واژه‌گزینی و ترجمه انتخاب می‌شوند، همه واژگان ثبت شده باید از این فرایند عبور کرده باشند و اگر واژه‌ای پیدا شود که در مجموعه مرجع وجود ندارد، قطعاً این واژه دارای غلط‌املائی یا حداقل ناسازگاری رسم‌الخطی است. شرط پذیرش پایگاه وردنت فاوا این است که همه واژگان باید در آزمون املا قبول شوند.

۲-۳-۱-۲- آزمون آوا

از آنجایی که تلفظ واژه به صورت الفبای استاندارد IPA ثبت شده است، می‌توان با استفاده از نگاشت این الفبا به الفبای فارسی صورت‌املائی واژه را بازسازی کرد. آزمون تلفظ این کار را انجام می‌دهد و واژگانی را که تلفظ آن‌ها با صورت‌املائی مطابق ندارد، استخراج می‌نماید.

۲-۳-۱-۳- آزمون عضویت در مترادف

در آزمون عضویت بررسی و تایید می‌شود که هر واژه باید حداقل عضو یک مترادف باشد.

۲-۳-۱-۴- ارزیابی دستی واژگان

در این نوع آزمون فهرست واژگان موجود در پایگاه وردنت فاوا در اختیار آزمونگر خبره قرار می‌گیرد. این آزمونگر تعدادی از واژه‌ها را انتخاب کرده و درستی مشخصات آن‌ها را بررسی می‌کند.


۲-۳-۲- آزمون مترادف

مهم‌ترین مفهوم بکار رفته در وردنت مترادف است که نشان‌دهنده مجموعه‌ای از کلمات یا عبارات هم‌معنی است. همچنین رابطه‌های بین این مجموعه مترادف‌ها مانند رابطه is-a (خاص/عام) از ویژگی‌های مهم وردنت می‌باشد. برای آزمون درستی مترادف‌ها و رابطه بین آن‌ها می‌توان فرایندهایی را بر مبنای وجود یک پیکره غنی متنی پیشنهاد نمود. همچنین یک فرهنگ لغت جامع می‌تواند برای ارزیابی اولیه بسیار موثر و مفید باشد. برای آزمون صحت و دقت مترادف‌ها، مترادف بودن واژگان عضو یک مترادف و همچنین درستی تعریف مترادف را بررسی می‌کنیم. در این راستا از منابع واژگانی مانند فرهنگ واژگان و پیکره متنی استفاده می‌نماییم.


۲-۳-۳- ارزیابی صحت و دقت رابطه‌های معنایی

برای ارزیابی درستی رابطه‌های معنایی شواهد و مراجع ایجاد آن رابطه را بررسی می‌کنیم. پیکره متنی و واژه‌نامه‌ها منابع اصلی برای ارزیابی روابط معنایی هستند. در فرهنگ لغت‌های معمول برای توصیف لوح فشرده از این عبارت استفاده می‌شود: "نوعی وسیله ذخیره‌سازی اطلاعات". در اینجا کلمه نوعی نشان‌دهنده یک رابطه خاص/عام بین لوح فشرده و وسیله ذخیره‌سازی اطلاعات است. با استفاده از این الگوها می‌توان رابطه فوق را بین کلمات یک مترادف با مترادف دیگر جستجو کرد.

۲-۳-۴- آزمون بازنمایی

	عنوان گزارش: گزارش آزمون پروژه وردنت فارسی در حوزه فاوا		پژوهشکده
	وضعیت گزارش: نهایی	کد گزارش: ITF.ITP.TCH.۱۳۹۹۱۳۲۰.V.۰۱	فناوری اطلاعات

برای ارزیابی کیفیت از آزمون بازنمایی استفاده می‌کنیم. مهم‌ترین مسئله در این آزمون تهیه سند شاهد است. این سند را به روش‌های مختلفی می‌توان ساخت. آسان‌ترین راه انتخاب تعدادی از منابع واژگانی قابل استفاده برای وردنت و عدم استفاده از آنها در ساخت وردنت می‌باشد. بدین ترتیب پس از ساخت وردنت، می‌توان محتوای آن را نسبت به این منابع سنجید.

	عنوان گزارش: گزارش آزمون پروژه وردنت فارسی در حوزه فاوا		پژوهشکده
	وضعیت گزارش: نهایی	کد گزارش: ITF.ITP.TCH.AMIRKOSAFAT.V.01	فناوری اطلاعات

نظارت تغییر یافت به نحوی که واژه های انگلیسی معادل یک واژه فارسی نیز در همان صفحه نمایش داده شود. مشکلات جزئی اندکی نیز مرتبط با نحوه نمایش اطلاعات در نرم افزار مرورگر رفع گردید.

۳-۳- آزمون صحت و اعتبار دامنه

دامنه های وردنت فاوا بر اساس تقسیم بندی ACM انتخاب شده اند که در آن همه سرفصل های اصلی پیشنهادی ACM پوشش داده شده اند. حوزه پوشش همه دامنه ها یکی نیست بنابراین نمی توان انتظار داشت تعداد واژگان همه منابع با هم مشابه باشد. با این حال تعداد واژگان هر دامنه با وسعت آن متناسب است.

۳-۴- آزمون صحت و اعتبار منابع واژگانی

فرهنگستان زبان و ادب فارسی، شرکت میکروسافت، امپریال کالج لندن ناشران منابع واژگانی اصلی هستند که همه از مراجع علمی معتبر به شمار می روند. کتب استفاده شده برای پیکره متنی نیز از ناشران شناخته شده و معتبر مانند Wiley, Prentice Hall, Mc Millan, MIT Press, Oxford University Press انتخاب شده اند.

واژگان مصوب فرهنگستان مرجع اصلی معادل های فارسی در همه حوزه هاست که به صورت چاپی و الکترونیکی در سطح وسیعی در کشور مورد استفاده قرار می گیرد. فرهنگ میکروسافت هم با شمارگان بالای ۳۰۰۰ جلد بارها در کشورمان تجدید چاپ شده است. سایت foldoc از پر بازدیدترین سایت های اینترنتی است. کتب استفاده شده برای پیکره حوزه وسیعی از مباحث فاوا را پوشش داده و از کتب مرجع در موضوع خود محسوب می شوند.


تلاش نمودیم تا همه واژگان این منابع را پوشش دهیم اما در مقاطع مختلف برخی از واژگان را به دلایل تفاوت فرهنگی یا غیر فاوا بودن آنها حذف نمودیم. با این وجود، همه واژگان تخصصی فاوای منتشر شده در دفترهای هفتم و هشتم فرهنگستان (جمعاً نزدیک به ۴۰۰ واژه) در وردنت فاوا ثبت شده است. همه واژگان فرهنگ میکروسافت و تقریباً همه واژگان foldoc پوشش داده شده اند. واژگان ذکر شده در منابع فوق به n-gram های پر بسامد پیکره متنی تعلق دارند و هر کدام از آنها حداقل ۱۰ بار در پیکره آمده اند.

۳-۵- آزمون های صحت و اعتبار واژه

تعریف هر واژه یا از فرهنگ های فوق بدست آمده و یا از ویکی پدیا استخراج شده است. با توجه به معتبر بودن منابع فوق تعریف واژه مورد وثوق صاحب نظران است.

۳-۵-۱- تعلق به حوزه فاوا

در وردنت فاوا، حوزه فاوا را به دامنه های متعددی تقسیم کرده ایم. بر همین اساس تنها واژگانی را که حداقل به یکی از این دامنه ها تعلق داشته اند را به عنوان واژه فاوا انتخاب نموده ایم. لازم به ذکر است که تعلق هر واژه به یک دامنه مرتبط با فاوا به صورت دستی و توسط متخصص فاوا انجام شده است که در

	عنوان گزارش: گزارش آزمون پروژه وردنت فارسی در حوزه فاوا		پژوهشکده
	وضعیت گزارش: نهایی	کد گزارش: ITF.ITP.TCH.AMIRKAZEM.V.01	فناوری اطلاعات

نتیجه آن شمار زیادی از واژگان (حدود ۵۰۰) واژه غیر مرتبط شناخته شد و از مجموعه دادگان وردنت حذف گردید.

۳-۵-۲- املا

صورت املائی واژه در وردنت انگلیسی به همان شکل ثبت شده در منبع مورد استفاده قرار گرفته است. واژگان زبان فارسی نیز مطابق قواعد املائی بروز فارسی نگارش شده‌اند.

۳-۵-۳- آوا

آوای واژگان فارسی بر اساس آوای مشهور و رایج آنها در محافل دانشگاهی ثبت شده است. آوای واژگان عاریتی نیز بر اساس آوای مرجع اصلی ثبت شده است. در وردنت فاوا برای هر واژه تنها یک آوا تعریف شده است. تعداد بسیار کمی واژه ممکن است دارای بیش از یک آوا باشند با این حال ما به چنین واژه ای برنخوریم.


۳-۵-۴- افزونگی

در وردنت فاوا صورت املائی هر واژه منحصر به فرد است. جدول‌های پایگاه داده طوری تعریف شده‌اند که از ورود واژه تکراری جلوگیری می‌شود. با این حال تنوع حروف فارسی با شکل مشابه و ب کدهای متفاوت سبب ایجاد افزونگی در جدول واژگان می‌شود. مثلاً حرف «ی» حداقل با سه کد ۱۶۰۹، ۱۶۱۰ و ۱۷۴۰ در جدول یونیکد تعریف شده است. سیستم عامل‌های مختلف در تخصیص کد حرف «ی» به صفحه کلید متفاوت عمل می‌کنند. همین عمل باعث می‌شود که واژگانی که صورت املائی مشابهی دارند اما این حرف با دو کد مختلف در آنها نوشته شده است، متفاوت قلمداد شوند. حرف «ک» هم وضعیت مشابهی دارد. در یک بررسی دریافتیم که ۸۳ واژه فارسی حداقل با دو شکل ثبت شده‌اند و بنابراین زاید هستند. برای غلبه بر این مشکل همه کدهای مختلف «ی» و «ک» را به یک کد تبدیل کردیم و این افزونگی برطرف شد.

۳-۶- آزمون صحت و اعتبار مترادف‌ها

در وردنت فاوا مترادف‌ها در زبان انگلیسی تعریف و به زبان فارسی ترجمه شده‌اند. ملاک قرار گرفتن واژگان در یک مترادف شباهت معنایی آنهاست. برای کشف این شباهت از روش‌های نیمه اتوماتیک و دستی استفاده شده است. در روش‌های نیمه اتوماتیک به کمک تحلیل‌های متنی تعریف واژگان، ابتدا دو یا چند واژه ای که تعریف آنها دارای تعداد قابل توجهی کلمه یا اصطلاح مشابه بود، انتخاب می‌شوند. سپس یک کارشناس خبره بری تجمیع این واژگان در یک مترادف تصمیم می‌گیرد.

۳-۶-۱- آزمون عضویت واژه در مترادف

	عنوان گزارش: گزارش آزمون پروژه وردنت فارسی در حوزه فاوا		پژوهشکده
	وضعیت گزارش: نهایی	کد گزارش: ITF.ITP.TCH.AMIRKAZEM.V.01	فناوری اطلاعات

این آزمون به صورت خودکار انجام می‌گیرد و هدف آن تعیین عضویت واژگان در ترادف‌هاست. با انجام این آزمون دریافتیم که در وردنت فاوا هر واژه حداقل در یک ترادف عضویت دارد و هر ترادف دارای حداقل یک واژه است.

۳-۶-۲ - آزمون شباهت معنایی واژگان

واژگان عضو یک ترادف باید هم معنی باشند. در وردنت انگلیسی فاوا واژگان بر اساس شباهت معنایی یا شواهد ساختاری در یک ترادف قرار گرفته‌اند. درستی شباهت معنایی واژگان ترادف‌ها را به صورت دستی بررسی می‌کنیم. در وردنت انگلیسی همه واژگان عضو ترادف دارای شباهت معنایی هستند. در وردنت فارسی فاوا ترجمه واژگان و اصطلاحاتی که دارای مخفف‌های یکسان هستند، با همدیگر ادغام شده است. این ادغام سبب ایجاد خطای عدم شباهت معنایی واژگان فارسی می‌گردد. در زیر مثال روشنی از این نوع خطا را آورده‌ایم:

بازیابی مبتنی بر محتوی|نرخ بیت ثابت → CBR|constant bit rate


بازیابی مبتنی بر محتوی|نرخ بیت ثابت → CBR|content-based retrieval

با بررسی که روی مخفف‌ها انجام دادیم، دریافتیم که نزدیک به ۵۰۰۰ ترادف امکان بروز چنین خطایی را دارند. این مجموعه از ترادف‌ها را به صورت دستی بررسی کرده و واژگان ناسازگار آنها را حذف کردیم. یکی دیگر از مشکلات ترجمه واژگان در وردنت فارسی، ترجمه اسامی خاص بود. مثلاً واژه rendezvous به صورت تحت‌اللفظی به پاتوق ترجمه شده است. در حالی که در حوزه فاوا نام یک زبان پرس و جو در پایگاه داده است. در مرحله بازبینی این واژه را به صورت «زبان rendezvous» ثبت کردیم. همچنین حدود ۱۵۰۰ مورد از واژگانی را که به صورت مشابه ترجمه شده بودند، کشف و اصلاح کردیم.

۳-۶-۳ - آزمون تعریف ترادف

تعریف ترادف‌های انگلیسی را از منبع واژگان آن ترادف استخراج کردیم. بخشی از تعاریف را از واژه نامه‌ها و برخی را به صورت خودکار از ویکی‌پدیا استخراج کردیم. در ادامه در هنگام بررسی دریافتیم که متن ثبت شده برای تعریف برخی از ترادف‌ها املا تعریف ترادف نیست و به جنبه دیگری مانند کاربرد ترادف می‌پردازد. علاوه بر این برخی تعاریف املا مفهوم دیگری را بیان می‌کردند. در این فاز حدود ۳۰۰۰ ترادف دارای تعریف نادرست، مبهم یا ناقص را پیدا کرده و تعریف آنها را اصلاح نمودیم. برای ترجمه تعاریف، آنها را بین گروهی از دانشجویان کارشناسی ارشد و دکتری مهندسی کامپیوتر تقسیم کردیم و سپس ترجمه‌ها را جمع‌آوری و تصحیح کردیم. از این میان تاکنون نزدیک به ۲۰ هزار تعریف را به فارسی ترجمه کرده و پس از تایید در پایگاه وردنت ثبت کرده‌ایم.

۳-۶-۴ - مقوله دستوری

	عنوان گزارش: گزارش آزمون پروژه وردنت فارسی در حوزه فاوا		پژوهشکده
	وضعیت گزارش: نهایی	کد گزارش: ITF.ITP.TCH.۱۳۹۹.۰۷.۰۱	فناوری اطلاعات


در وردنت فاوا مقوله دستوری در سطح ترادف تعریف می‌شود و مقادیر معتبر آن اسم، فعل، صفت و قید می‌باشد. برای حدود ده هزار ترادف مقوله دستوری از فرهنگ لغت مایکروسافت استفاده کردیم و مقوله بقیه ترادف‌ها را به صورت دستی تعیین کردیم. در فرایند آزمون یک‌بار دیگر مقوله دستوری ترادف را با تعریف آن بررسی کردیم و دریافتیم که حدود ده ترادف که برچسب اسم داشتند، املا به مقوله صفت تعلق دارند. همچنین هفت ترادف اسمی به مقوله فعل و یک ترادف اسمی هم به مقوله قید متعلق بودند.

۳-۶-۵- دامنه اصلی و فرعی

در وردنت فاوا هر ترادف دارای یک برچسب دامنه سطح اول و چند برچسب دامنه سطح دوم می‌باشد. هر دو برچسب دامنه به صورت دستی تعیین شده‌اند. در این فاز برچسب دامنه همه ترادف‌ها را به صورت دستی بررسی و اصلاح کردیم.

۳-۷- صحت و اعتبار رابطه های معنایی

رابطه‌هایی معنایی عمدتاً به کمک برنامه های پردازش متن از پیکره متنی یا ویکی‌پدیا استخراج شده و سپس به صورت دستی تصفیه شده‌اند. در مرحله اول بیش از ۱۰۰ هزار رابطه معنایی به صورت اتوماتیک و نیمه اتوماتیک از پیکره متنی و واژه نامه‌ها استخراج کردیم. سپس در مرحله دوم با بازبینی انسانی به وسیله تیمی از دانشجویان کارشناسی و کارشناسی ارشد، از این میان ۳۰ هزار رابطه را تایید کردیم. در مرحله سوم تاکنون ۲۲ هزار رابطه را توسط تیم متخصص فاوا شامل اعضای هیئت علمی بررسی و در پایگاه وردنت ثبت نمودیم. در جدول ۳-۱ آمار روابط معنایی استخراج شده، ثبت شده و تایید شده آمده است.


	عنوان گزارش: گزارش آزمون پروژه وردنت فارسی در حوزه فاوا		پژوهشکده
	وضعیت گزارش: نهایی	کد گزارش: ITF.ITP.TCH.AMIRKAZEM.V.01	فناوری اطلاعات

جدول ۳-۱ (آمار روابط معنایی وردنت فاوا

کد	نام رابطه	موجود	ثبت شده	پیکره	وب	تایید انسانی	تایید نشده
۱	hypernym		۱۹۳۰۰			۱۹۳۰۰	۰
۲	hyponym		۱۹۳۰۰			۱۹۳۰۰	۰
۳	instance hypernym		۳۷۱			۳۷۱	۰
۴	instance hyponym		۳۷۱			۳۷۱	۰
۱۱	part holonym		۲۶			۲۶	۰
۱۲	part meronym		۲۶			۲۶	۰
۱۳	member holonym		۵۱			۵۱	۰
۱۴	member meronym		۵۱			۵۱	۰
۱۵	substance holonym		۳			۳	۰
۱۶	substance meronym		۳			۳	۰
۲۱	entailment		۹			۹	۰
۲۳	causes		۱۳			۱۳	۰
۳۰	antonym		۵۴۶			۵۴۶	۰
۴۰	similar		۶			۶	۰
۵۰	see_also		۳۸۷۰			۳۸۷۰	۰
۶۰	attribute		۴۶			۴۶	۰
۹۱	domain category	۵۷۴۸۶	۵۷۴۸۶			۳۰۰۰۰	۰
۹۲	domain member category	۵۷۴۸۶	۵۷۴۸۶			۳۰۰۰۰	۰
۲۰۱	comes after		۱۷			۱۷	۰
۲۰۲	has_attribute		۲۵۶			۲۵۶	۰
۲۰۳	has_function		۲۲۸			۲۲۸	۰
۲۰۴	is_used_for		۱۱۰			۱۱۰	۰
۲۰۷	uses		۱۱۰			۱۱۰	۰
۲۲۱	subject		۲۹۱			۲۹۱	۰
۲۲۲	adv complement		۱			۱	۰
۲۲۳	object		۲۷۶			۲۷۶	۰
۲۲۴	Indirect object		۹۳			۹۳	۰
۲۲۵	modifies		۳۳۶			۳۳۶	۰
۲۲۶	modified_by		۳۳۶			۳۳۶	۰

۳-۸- نگاشت وردنت انگلیسی به فارسی

در وردنت فاوا هر مترادف انگلیسی دقیقاً به یک مترادف فارسی ترجمه شده است و متقابلاً هر مترادف فارسی هم دقیقاً به یک مترادف انگلیسی نگاشت یافته است. رابطه های معنایی بین مترادفها در دو وردنت

	عنوان گزارش: گزارش آزمون پروژه وردنت فارسی در حوزه فاوا		پژوهشکده
	وضعیت گزارش: نهایی	کد گزارش: ITF.ITP.TCH.AMIRKHAZEM.V.01	فناوری اطلاعات

دقیقاً مثل هم است. به این دلیل نگاشت دو وردنت به همدیگر یک به یک و پوشاست. در پایگاه داده وردنت شناسه مترادف^۱ در هر دو زبان یکی است.

۳-۹- آزمون کیفی محتوا

در این دسته آزمون‌ها کیفیت دادگان ثبت شده در وردنت ارزیابی می‌شود. مهم‌ترین این آزمون‌ها آزمون بازنمایی، آزمون افزونگی مترادف و آزمون دور در روابط شمول معنایی هستند.

۳-۹-۱- آزمون بازنمایی


برای آزمون بازنمایی مجموعه واژگان TechTerms را در قالب سند شاهد استفاده کردیم. ۵۰۰ واژه از این مجموعه واژگان را که دارای ضریب^۲ ۵ یا کمتر بودند در پایگاه داده جستجو کردیم. از این میان ۲۰ واژه در وردنت فاوا موجود نبودند. پس از بررسی دریافتیم که تنها ۷ واژه password, numlock, iOS, address, bar, refresh, kindle, iPod در وردنت وجود ندارد و بقیه به صورت‌های دیگر املایی ثبت شده‌اند. مثلاً keyboard به صورت keyborads ثبت شده بود. این هفت واژه را در وردنت فاوا ثبت کردیم. به کمک تعاریف موجود در TechTerms تعریف ۵ مترادف را اصلاح و بهسازی کردیم. قابل ذکر است که علاوه بر روابط ثبت شده در وردنت روابط معنایی قابل توجهی نیز استخراج گردیده است که می‌توان در آینده پس از پالایش انسانی به دادگان وردنت فاوا اضافه و آن را غنی‌تر نمود.

۳-۹-۲- آزمون افزونگی مترادف

با توجه به تولید خودکار بخش اعظم وردنت فاوا، یکی از چالش‌های به وجود آمده افزونگی مترادف‌ها می‌باشد. روش خودکار پیشنهادی ما برای تشکیل مترادف‌ها پردازش متن تعریف آنها بود. به این صورت که واژگانی که تعاریف آنها شباهت بیشتری به هم داشت در یک مترادف قرار می‌گرفت. مشکل این روش این است که بسیاری از تعاریف با اینکه از لحاظ مفهومی کاملاً یکسان هستند، از لحاظ واژگان و نیز طول جمله تعریف با همدیگر تفاوت دارند و به همین جهت الگوریتم پیشنهادی آنها را در یک مترادف قرار نمی‌دهد. متأسفانه تعداد مترادف‌های تکراری در وردنت فاوا قابل توجه بود که برای رفع این مشکل ناچار شدیم که مترادف‌های افزونه را به صورت دستی پیدا، حذف و در مواردی با هم ادغام کنیم. در این فاز بیش از ده هزار مترادف را بررسی و مترادف‌های تکراری را در همدیگر ادغام نمودیم. لازم به ذکر است که در فرایند ادغام مترادف‌ها هیچ واژه‌ای از دست نمی‌رود. مزیت دیگر ادغام این مترادف‌ها این است که اگر یک مترادف به صورت جداگانه دارای رابطه معنایی با مترادف دیگری نباشد، پس از ادغام شانس بسیار بیشتری برای اتصال به شبکه معنایی دارد.

^۱ SynsetID

^۲ Tech Factor

	عنوان گزارش: گزارش آزمون پروژه وردنت فارسی در حوزه فاوا		پژوهشکده
	وضعیت گزارش: نهایی	کد گزارش: ITF.ITP.TCH.۱۳۹۳.۱۰.۷۰۰	فناوری اطلاعات


۳-۱۰- ارزیابی عوامل و ملاک‌ها

بر اساس آزمون‌های انجام شده، عوامل و ملاک‌ها ارزیابی می‌شوند. نتایج ارزیابی عوامل و ملاک‌ها در جدول زیر خلاصه شده است. به جز رابطه شمول معنایی و ترجمه تعاریف، در سایر عوامل سطوح مطلوب تأمین شده است.

جدول (۲-۳) ارزیابی عوامل، ملاک‌ها و نشانگرها


مطلوبیت؟	سطح مطلوب	ملاک / نشانگر	عامل
۱۰۰	همه منابع باید معتبر باشند	اعتبار	منابع واژگانی
۱۰۰	همه منابع باید جزو منابع مشهور و پرکاربرد باشند	پرکاربرد بودن	
۱۰۰	همه منابع باید غنای واژگانی کافی داشته باشند	غنای واژگانی	
۱۰۰	واژگان منتخب باید از جمله واژگان پرسامد فاوا باشند	پر بسامد بودن	مجموعه واژگان
۱۰۰	همه زیر دامنه‌های تعریف شده باید واژگانی داشته باشند	پوشش زیر دامنه‌ها	
۱۰۰	توزیع واژگان در زیر دامنه‌ها باید همگن باشد.	کفایت از نظر تنوع	
۱۰۰	کل واژگان باید حداقل ۳۰۰۰۰ واژه باشد	کفایت از نظر تعداد	
۱۰۰	تعریف هر واژه باید جامع و مانع باشد	درستی تعریف	واژه
۱۰۰	املاي همه واژه‌ها به صورت صحیح ثبت شود	درستی املا	
۱۰۰	صورت آوایی همه واژه‌ها باید به صورت صحیح ثبت شود	درستی تلفظ	
۱۰۰	تعاریف اعضای مترادف‌ها باید مشابه باشد	شباهت معنایی اعضا	ترادف
۱۰۰	تعریف هر مترادف باید جامع و مانع باشد	صحت تعریف	
۱۰۰	همه رابطه‌های تعریف شده بین مترادف‌ها باید صحیح باشند	صحت رابطه	رابطه معنایی
۱۰۰	منبع مورد استفاده برای هر رابطه معتبر باشد..	منبع استخراج	
۱۰۰	حداقل پنج نوع رابطه معنایی باید تعریف شود.	انواع روابط	پوشش رابطه معنایی
۷۵	همه واژگان مقوله‌های اسم و فعل دارای مترادف پدر باشند.	شمول معنایی	
۱۰۰	حداقل ۲۰۰ رابطه اشتقاق بین واژگان فارسی استخراج شود.	پوشش اشتقاق	پوشش
۱۰۰	حداقل ساختار ظرفیتی ۲۰۰ فعل قید شود	ساختار ظرفیتی افعال	روابط ساختاری
۱۰۰	حداقل روابط آرگومانی ۲۰۰ فعل قید شود	روابط آرگومانی افعال	
۱۰۰	همه مترادف‌های فارسی به انگلیسی و بالعکس نگاشت یابند	وجود نگاشت	نگاشت وردنت
۱۰۰	هر مترادف به مترادف هم معنی نگاشت یابد	صحت نگاشت واژگان	
۷۰	تعاریف انگلیسی به معادل مناسب فارسی نگاشت یابند.	صحت نگاشت تعاریف	

	عنوان گزارش: گزارش آزمون پروژه وردنت فارسی در حوزه فاوا		پژوهشکده فناوری اطلاعات
	وضعیت گزارش: نهایی	کد گزارش: ITF.ITP.TCH.۱۳۹۹۱۳۲۰۰۰.V.۰۲	
شد.	دور و پل در رابطه های معنایی وردنت فاوا پیدا	آزمون دور و پل در روابط	
نشد.	با بررسی دستی حداقل ۵۰۰۰ مترادف در هم ادغام	آزمون افزونگی مترادف	
شد. حداکثر حدود همین تعداد باقی مانده است.	مراجعه به جدول مربوطه	ارزیابی عوامل و ملاکها	ارزیابی عوامل و ملاکها
			۸


	عنوان گزارش: گزارش آزمون پروژه وردنت فارسی در حوزه فاوا		پژوهشکده
	وضعیت گزارش: نهایی	کد گزارش: ITF.ITP.TCH.۱۳۹۳۱۳۰۰.V.۰۲	فناوری اطلاعات

۴- جمع بندی


آزمون پروژه وردنت فاوا از جنبه‌های کارکردی، غیر کارکردی و محتوایی انجام می‌گیرد. در این سند انواع آزمون‌های ممکن در جنبه‌های فوق را معرفی کرده و نتایج آنها را ارائه کردیم.

	عنوان گزارش: گزارش آزمون پروژه وردنت فارسی در حوزه فاوا		پژوهشکده
	وضعیت گزارش: نهایی	کد گزارش: ITF.ITP.TCH.۱۳۹۹۱۳۲۰.V.۰۲	فناوری اطلاعات

English	فارسی
Test adaptor	آزمونگر
Test case	مورد آزمون
Test designer	طراح آزمون
Test driver	راه انداز آزمون
Test manager	مدیر آزمون
Test plan	طرح آزمون
Test scope	گستره آزمون
Thesaurus	اصطلاح نامه
Unit testing	آزمون واحد
Usability testing	آزمون قابلیت استفاده
User Interface	واسط کاربر
Volume testing	آزمون حجم
White-box Testing	آزمون جعبه سفید
Word	واژه

	عنوان گزارش: گزارش آزمون پروژه وردنت فارسی در حوزه فاوا		پژوهشکده فناوری اطلاعات
	وضعیت گزارش: نهایی	کد گزارش: ITF.ITP.TCH.۱۳۹۳۱۳۲۰.V.۰۲	

فارسی	English
قواعد انجمنی	Association Rules
کلاس	Class
گستره آزمون	Test scope
متدولوژی RUP	RUP Methodology
مدیر آزمون	Test manager
مدیریت دانش	Knowledge Management
موتور استنتاج	Inference Engine
مورد آزمون	Test case
مؤلفه	Component
نقش	Role
نیازمندی‌های غیر کارکردی	Non-functional Requirements
نیازمندی‌های کارکردی	Functional Requirements
هستان شناسی	Ontology
واژه	Word
واسط کاربر	User Interface

	عنوان گزارش: گزارش آزمون پروژه وردنت فارسی در حوزه فاوا		پژوهشکده فناوری اطلاعات
	وضعیت گزارش: نهایی	کد گزارش: ITF.ITP.TCH.۱۳۹۳.۱۳.۱. V.۲	

Abstract

This document presents test results for our WordNet in ICT domain.



Information Technology Faculty

Information Technology Platform Group

Technical Report

Test Reports

Project Name: Persian WordNet for ICT Domain

Project code: ۸۹۳۲۴۱۶

Project Director	Muharram Mansoorizadeh
Author(s)	M. Mansoorizadeh, N. Nassiri and M. Dadrass
Document Code	ITF.ITP.TCH.۸۹۳۲۴۱۶.۴۰.V.۲
Preparing Date	۹۶.۱۱.۱۱
Status/Version	Final/۲.۰