

استخراج پیکره موازی از اسناد قابل مقایسه برای بهبود کیفیت ترجمه در سامانه‌های ترجمه ماشینی

زینب رحیمی، محمدحسین ثمنی و شهرام خدیوی

دانشکده مهندسی کامپیوتر و فناوری اطلاعات، دانشگاه صنعتی امیرکبیر، تهران، ایران

چکیده

امروزه با گسترش وسایل ارتباط عمومی و به خصوص شبکه جهانی اینترنت، نیاز به عملیات ترجمه خودکار به صورت چشم‌گیری افزایش یافته است. یکی از مطرح‌ترین روش‌های ترجمه ماشینی، روش آماری است. پارامترهای سامانه^۱ ترجمه ماشینی آماری با استفاده از مجموعه بزرگی از دادگان آموزشی (پیکره موازی دوزبانه) تخمین زده می‌شود؛ اما در برخی زبان‌ها، هنوز مسأله نیاز پایه‌ای سامانه ترجمه ماشینی آماری یعنی پیکره‌های متنی بزرگ موازی برطرف نشده است. برای رفع این مشکل روشی پیشنهادی جهت بهبود کیفیت پیکره‌های مستخرج از اسناد قابل مقایسه و در نتیجه بهبود کیفیت سامانه ترجمه ماشینی ارائه شده است. از آنجایی که در اکثر متون قابل مقایسه داده‌های موازی نه به صورت جمله، بلکه به صورت قطعات زیرجمله‌ای ظاهر می‌شوند، روش پیشنهادی سعی در استخراج قطعات موازی به صورت بلوک با استفاده از مجموعه‌ای از ویژگی‌ها دارد که این ویژگی‌ها عبارت‌اند از طول عبارت، امتیاز شباهت لگاریتمی، شیب مسیر ترازبندی در بلوک، پراکندگی شیب قطعات تشکیل‌دهنده بلوک، مربعی بودن بلوک و درصد حضور کلمات هم‌ترجمه در بلوک. طبق ارزیابی‌های انجام‌شده روش پیشنهادی کارایی مناسبی دارد؛ و علاوه بر اینکه از نظر دقت و بازخوانی از روش‌های موجود استخراج قطعه پیشی گرفته است، دادگان مستخرج از اجرای این روش روی، بخشی از پیکره قابل مقایسه موجود، کارایی سامانه ترجمه ماشینی پایه را برای دادگان آزمون مختلف از ۰/۳۳ تا ۱/۴ واحد بلو افزایش داده است.

واژگان کلیدی: پیکره قابل مقایسه^۲، استخراج قطعات موازی^۳، پیکره موازی^۴، ترجمه ماشینی^۵

۱- مقدمه

با توجه به حجم روزافزون اطلاعات و اهمیت پردازش خودکار گفتار و اسناد توسط رایانه، می‌توان به اهمیت مدل‌سازی زبان و انجام عملیات ترجمه توسط ماشین پی برد. در زمینه ترجمه ماشینی دو رویکرد اصلی مبتنی بر قانون و روش آماری وجود دارند که روش مورد نظر ما در این پژوهش روش آماری است؛ یعنی استفاده از مدل‌سازی آماری برای ترجمه یک متن از یک زبان به زبانی دیگر. رویکرد آماری ترجمه ماشینی یکی از روش‌های اصلی در میان روش‌های مبتنی بر پیکره نیز هست؛ چون سامانه

طراحی‌شده، تمام اطلاعات مورد نیاز خود را از یک پیکره موازی دوزبانه متشکل از مجموعه بسیار بزرگی از جملات هم‌ترجمه استخراج می‌کند.

مبنای عملکرد یک مترجم آماری، نظریه تصمیم آماری است. نظریه تصمیم آماری روش شناخته‌شده‌ای برای ساخت یک سامانه تصمیم‌گیری مرکب، از چندین منبع اطلاعاتی موجود، با هدف حداقل‌سازی خطای تصمیم‌گیری است. در ترجمه ماشینی آماری، بر آنیم تا احتمال ترجمه $\Pr(f_1^j | e_1^j)$ را به نحوی مدل کنیم که رابطه بین رشته زبان مبدأ f_1^j و رشته زبان مقصد e_1^j را بیان کند. شکل (۱-۱) روند کلی ترجمه را در سامانه آماری نشان می‌دهد که در آن I و J به ترتیب تعداد کلمات رشته زبان مبدأ و مقصد هستند. پارامترهای چنین سامانه‌ای با استفاده از مجموعه بزرگی از دادگان آموزشی (پیکره موازی دوزبانه) تخمین زده می‌شود.

¹ System

² Comparable corpus

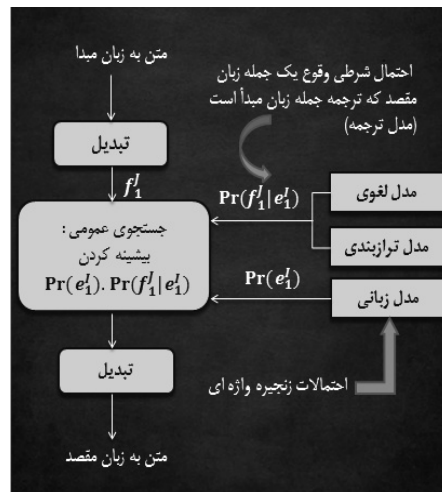
³ Parallel fragments extraction

⁴ Parallel corpora

⁵ Machine translation

این روش مدل‌هایی را با داده‌های آموزشی بسیار زیاد آموزش می‌دهد و بنابراین قدرت محاسباتی زیادی احتیاج دارد. پارامترهای آماری مورد نیاز برای مدل آماری به صورت خودکار از توزیع کلمات در دادگان متنی بزرگ استخراج و از جفت جملات موازی ورودی و خروجی یاد گرفته می‌شوند. در واقع یک مدل آماری، با توجه به ویژگی‌های متفاوت هم‌وقوعی و غیرهم‌وقوعی یک کلمه خاص با دیگر کلمات در زبان دیگر و با توجه به یک پیکره آموزشی دوزبانه یاد گرفته می‌شود.

در سال‌های اخیر برای برخی زبان‌ها از روش آماری استفاده شده است؛ اما در برخی زبان‌ها هنوز مسأله نیاز پایه‌ای سامانه ترجمه ماشینی آماری یعنی پیکره‌های متنی بزرگ موازی برطرف نشده است.



(شکل ۱-۱): روند کلی ترجمه در سامانه ترجمه ماشینی آماری

روش آماری به کمک سه منبع، داده‌های آموزشی خود را به دست می‌آورد. این سه منبع پیکره‌های متنی موازی، غیرموازی و قابل مقایسه بوده که جزء تفکیک‌ناپذیری از سامانه‌های ترجمه آماری پیکره‌های موازی هستند. پیکره‌های موازی منبعی ایده‌آل برای ترجمه هستند؛ اما به خاطر کمبود این متون از اسناد قابل مقایسه برای استخراج متون موازی استفاده می‌کنیم؛ حتی با وجود اینکه این متون برای زبان‌های خاصی موجود هستند، ممکن است در دامنه خاصی موجود نباشند و استفاده از متون موازی با زمینه متفاوت برای آموزش سامانه ترجمه، از کارایی سامانه ترجمه ماشینی به شدت خواهد کاست؛ پس به ناچار از متون غیرموازی و قابل مقایسه در کنار آنها استفاده می‌کنیم. از آنجایی که پیکره‌های موازی در آموزش سامانه‌های ترجمه

ماشینی نقش اساسی دارند، بر آنیم تا این پیکره‌ها را با استفاده از اسناد قابل مقایسه ایجاد کنیم.

منظور از پیکره‌های قابل مقایسه، اسنادی با موضوعات مشابه در دو زبان متفاوت است که با توجه به رشد و گسترش روزافزون وسایل ارتباط عمومی و اطلاع‌رسانی، به خصوص اینترنت، در حجم زیادی موجود هستند. به عنوان مثال، اخبار منتشرشده از وقایع مشترک در زبان‌های مختلف، نمونه‌ای از این نوع پیکره‌هاست. مزیت این گونه منابع در دسترس بودن آنهاست که موجب شده تا ایجاد و استفاده از پیکره‌های قابل مقایسه به عنوان زمینه ارزشمندی در بازیابی اطلاعات بین زبانی به خصوص در وب مطرح شود.

برای تهیه پیکره موازی دو رویکرد عمده روش دستی و خودکار وجود دارد. در این راستا ابتدا باید متونی که با نیازها مطابقت داشته باشند، را جمع‌آوری کرد؛ گاهی لازم است داده‌های اولیه جمع‌آوری شده و به صورت دستی وارد نوشتاری را به متون الکترونیکی تبدیل کرد. این کار بسیار زمان‌بر است و انجام آن به صورت دستی و یا با استفاده از نرم‌افزار OCR به صرفه نیست. OCR با پوشش کردن و تشخیص کلمات این کار را انجام می‌دهد. از سوی دیگر یک سری داده‌های از قبل منتشرشده، از قبیل متن‌های اداری و خبری موجود هستند و این گونه داده‌ها مانند صفحات وب، صفحات شخصی افراد و یا شرکت‌ها، فهرست سخنرانی‌ها در اینترنت در حجم زیاد قابل استفاده هستند. بدیهی است که استفاده از روش‌های خودکار برای انجام فرآیند ایجاد پیکره موازی یا قابل مقایسه و در گام بعدی یعنی استخراج دادگان موازی، مناسب‌ترین روش است.

پس از جمع‌آوری متون، برای یافتن جملات موازی (هم ترجمه) در آنها جستجو انجام می‌گیرد و جملات موازی متناظر به یکدیگر تراز می‌شوند. روش‌های مختلفی برای این ترازبندی وجود دارد. یک روش خوب باید سریع و دقیق باشد، به زبان وابسته نباشد و احتیاج به دانش خاصی در مورد دو زبان نداشته باشد. اگر سندها ترازبندی نشده باشند و ترازبندی در حد سند موجود نباشد، نیاز به توسعه روش‌هایی است که این ترازبندی را انجام دهند.

روش‌های ترازبندی در سطوح مختلفی مانند کلمه، قطعه، جمله، بند و سند (برای متون قابل مقایسه) انجام می‌شوند. در ترازبندی براساس کلمه، از مدل‌های آماری استفاده می‌شود. در ترازبندی در سطح قطعات، متن به قطعاتی تقسیم می‌شود که اگر با روش احتمالاتی، ترجمه

از دیگر فعالیت‌های پژوهشی مربوط به این زمینه کار رزینیک و/اسمیت در سال ۲۰۰۳ (رزینیک و اسمیت، ۲۰۰۳) است که سامانه آنها برای کشف جفت‌اسناد موازی در وب طراحی شده است.

اولین تلاش در زمینه کاوش پیکره‌های قابل مقایسه برای استخراج جملات موازی توسط زائو و وگل در سال ۲۰۰۲ انجام شد که از گسترش الگوریتم‌های طراحی شده برای ترازبندی جملات متون موازی استفاده کرده و با استفاده از برنامه‌نویسی پویا اسنادی را که فرض می‌شود، مشابه‌اند ترازبندی کردند. این روش‌ها تنها برای "اسناد موازی و نوفه‌ای" مناسب است؛ یعنی برای اسنادی که به‌نسبه مشابه هستند، هم در محتوا و هم ترتیب جملات قابل اجرا می‌باشد.

در کار (بارزیلا و الحداد، ۲۰۰۳) ترازبندی جملات با استفاده از ویژگی شباهت ترجمه انجام شده است. به این صورت که برای ترازکردن جملات، شباهت لغوی آنها را در نظر گرفته است. تحقیق گزارش شده توسط فونگ و چونگ (فونگ و چونگ، ۲۰۰۴) و (فونگ و وو، ۲۰۰۵) نیز برای اسناد غیرموازی عمل می‌کند. آنها نیز هر سند مبدأ را با چند سند مقصد متناظر کرده و تمام جفت‌جمله‌های ممکن را بررسی می‌کنند؛ اما فهرست ثابتی از جفت‌اسناد وجود ندارد. پس از یک دور استخراج جملات، فهرست با اسناد بیشتر کامل شده و سامانه به‌صورت تکراری اجرا می‌شود. مانتینیو و مارکو در پژوهش سال ۲۰۰۵ خود (مانتینیو و مارکو، ۲۰۰۵) استفاده از بیش‌ترین بی‌نظمی را برای انتخاب جملات هم‌تراز استفاده کردند؛ بدین شکل که جفت‌جملات را جدا از مضمون جمله تجزیه و تحلیل کرده و آنها را به‌عنوان موازی و غیر موازی طبقه‌بندی کردند. آنها هر سند مبدأ را با چندین سند مقصد مطابقت داده و همه جفت‌جملات ممکن را از هر جفت سند طبقه‌بندی کردند. این کار آنها را قادر می‌سازد که جملات از اسناد به‌نسبه غیرمشابه پیدا کنند و هر مقدار از تغییر ترتیب جملات قابل تحمل باشد.

در روش اکس ما (اکس ما، ۲۰۰۶) یک مدل رتبه‌بندی‌کننده آموزش می‌بیند که برای هر جمله در سند مبدأ یک جمله موازی در سند مقصد انتخاب می‌کند یا تهی را به آن نسبت می‌دهد. این فرمول‌سازی از مشکل عدم توازن طبقه‌بندی‌کننده دودویی جلوگیری می‌کند و در هر دو روش طبقه‌بندی‌کننده دودویی و رویکرد رتبه‌بندی‌کننده از طبقه‌بندی‌کننده بیش‌ترین بی‌نظمی استفاده می‌شود.

یکدیگر تشخیص داده شدند، جفت‌قطعه نامیده می‌شوند. این کار (ترازبندی زیر جملات (قطعات)، به جای ترازبندی جملات) که در اسناد موازی انجام می‌شود، حافظه مورد استفاده را کاهش می‌دهد و کیفیت فرآیند ترازبندی لغات را بهبود می‌بخشد. بندها با استفاده از شباهت لغوی یا خوشه‌بندی مبتنی بر محتوا ترازبندی و به‌طور کلی در ترازکردن اسناد نیز از کلمات مشترک آنها استفاده می‌شود.

گفتیم برای تهیه پیکره موازی دو رویکرد عمده دستی و خودکار وجود دارد؛ همچنین به معرفی پیکره‌های قابل مقایسه و حجم قابل توجه آن در وب پرداختیم. در اینجا سعی شده تا روشی برای بهبود کیفیت پیکره‌های مستخرج از اسناد قابل مقایسه و در نتیجه بهبود کیفیت سامانه ترجمه ماشینی ارائه شود. از آنجایی که در اکثر متون قابل مقایسه داده‌های موازی، نه به‌صورت جمله، بلکه به‌صورت قطعات زیرجمله‌ای ظاهر می‌شوند، روش پیشنهادی سعی در استخراج قطعات موازی به‌صورت بلوک با استفاده از مجموعه‌ای از ویژگی‌ها دارد.

در ادامه در بخش دوم به معرفی اهم کارهای انجام‌شده در این زمینه پرداخته می‌شود. در بخش سوم به معرفی روش پایه پیاده‌سازی شده پرداخته شده و روش جدید پیشنهادی معرفی و جزئیات آن بیان می‌شود. در بخش ارزیابی به ارزیابی روش‌های پیاده‌سازی شده و تحلیل نتایج پرداخته می‌شود و در نهایت جمع‌بندی و نتیجه‌گیری انجام می‌شود.

۲- کارهای انجام‌شده

برای استخراج پیکره موازی از اسناد موازی، روش‌های مشهوری مانند روش (گیل و چرچ، ۱۹۹۱) وجود دارند که در آن ترازبندی با طول جمله شباهت لغوی (مدل آی‌بی‌ام ۱) انجام می‌شود؛ روش (مور، ۲۰۰۲) که از تطابق کلمات و طول جملات برای ترازبندی و روش (ملاسد، ۱۹۹۹) که برای پیدا کردن ترازبندی‌های بین دو متن، از خصوصیات هندسی آنها استفاده می‌کنند؛ اما در زمینه پیکره‌های قابل مقایسه مطالعات محدودتر است. بخش عمده‌ای از کارهایی که با پیکره‌های قابل مقایسه انجام شده‌اند روی استخراج ترجمه کلمه متمرکز هستند. (فونگ و یی، ۱۹۹۸؛ دیاب و فیسنج، ۲۰۰۰؛ کوهن و نایت، ۲۰۰۴؛ گاوسیر و همکاران، ۲۰۰۴)

¹IBM 1

روش (تیلمن، ۲۰۰۹) در واقع توسعه‌ای برای روش ارائه‌شده توسط مانتینیو است که در آن برای جمله‌داده‌شده، یک طبقه‌بندی‌کنندهٔ بیش‌ترین آنتروپی به مجموعه بزرگی از ترجمه‌های نامزد اعمال می‌شود. یک الگوریتم جستجوی شعاعی برای این کار استفاده می‌شود که جملات مقصدی که خارج از شعاع^۱ می‌افتند هر چه سریع‌تر غیرموازی تشخیص داده شوند. در روش (عبدالرئوف و همکاران، ۲۰۰۹) آزمایش‌ها با استفاده از فیلتر جدید نرخ خطای ترجمهٔ امتیازی در کنار فیلترهای نرخ خطای ترجمه و نرخ خطای کلمه انجام شده است و در آن از یک سامانهٔ ترجمهٔ ماشینی آماری استفاده شده که از دادگان موازی کمی برای ترجمه سمت مبدأ استفاده می‌کند و متون مقصد در قسمت مدل زبانی سامانه استفاده می‌شوند. پس از آن از تکنیک‌های بازبایی اطلاعات و فیلترهای ساده برای ایجاد دادهٔ موازی از پیکرهٔ خبری قابل مقایسه استفاده می‌شود. روش (تیلمن و ژو، ۲۰۰۹) روش ساده‌ای است؛ به این صورت که مجموعهٔ بزرگی از جملات نامزد در سطح جمله به‌طور مستقیم امتیازدهی می‌شوند. در این روش استخراج در سطح جملات، وابسته به پیاده‌سازی کارآمدی از تابع امتیازدهی متقارن ساده است که منجر به افزایش سرعت محاسبه با عامل ۳۰ شده است.

در روش (لی و همکاران، ۲۰۱۰) یک مدل ترکیبی بدون ناظر معرفی می‌شود که ویژگی‌های آماری، لغوی، زبانی، مفهومی و زمانی را در یک چارچوب مینی بر حداکثرسازی انتظار معمولی را ترکیب می‌کند تا واژگان دوزبانه را از پیکرهٔ قابل مقایسه استخراج کند.

یکی از محدودیت‌های همهٔ این روش‌ها این است که آنها برای پیدا کردن جملات کامل طراحی شده‌اند. روش ما این است که چون جملات به‌طور کامل موازی در اسناد قابل مقایسه کم هستند، زیرجملات موازی را تشخیص دهیم. کارهای زیادی در این زمینه انجام نشده است. کار اصلی انجام‌شده در زمینهٔ استخراج قطعهٔ سامانهٔ طراحی شده توسط مانتینیو و مارکو در سال ۲۰۰۶ است (مانتینیو و مارکو، ۲۰۰۶) که با تجزیه و تحلیل جفت‌جملات به‌طور بالقوه مشابه، با استفاده از روشی الهام‌گرفته از پردازش سیگنال، تشخیص می‌دهد که کدام بخش از جملهٔ مبدأ به کدام بخش از جملهٔ مقصد مرتبط است.

کاردنگ و همکاران (دنگ و همکاران، ۲۰۰۶) نیز به قطعات زیرجمله می‌پردازد؛ اما آنها قطعات موازی را از

^۱beam

جفت‌جمله‌های موازی استخراج می‌کنند؛ درحالی که ما آنها را از جفت‌جملات قابل مقایسه و یا غیر موازی به‌دست می‌آوریم. در روش (هوایتارانا، ۲۰۱۰) از سه روش ترازبندی عبارات برای یافتن جفت عبارات استفاده می‌شود: الگوریتم استاندارد استخراج عبارات (ابزار موزس و با استفاده از مکاشفه)، استفاده تنها از ویژگی‌های نحوی که وابسته به مسیر ویتربی نیست و طبقه‌کنندهٔ دودویی با استفاده از روش بی‌نظمی بیشینه.

در زمینهٔ استخراج متون موازی از اسناد قابل مقایسه برای زبان فارسی کارهای شاخص زیادی انجام یا منتشر نشده است. یکی از محدود کارهای انجام‌شده، توسعهٔ پیکرهٔ تطبیقی UTPECC (هاشمی و همکاران، ۲۰۱۰) است که از تکنیک‌های بازبایی اطلاعات استفاده کرده و با تکیه بر تاریخ انتشار اخبار و شباهت محتوای اسناد عمل می‌کند. کار دیگر مرتبط (محمدی و قاسم آقایی، ۲۰۱۰) است که در این مقاله یک پیکرهٔ متنی موازی و ترازشده برای جفت زبان‌های فارسی و انگلیسی با کاوش در محتویات ویکی‌پدیا ارائه می‌شود. در این پژوهش روشی برای ترازبندی مبتنی بر لغت‌نامه دوزبانه در سطح جمله ارائه شده که از روشی بهبودیافته مبتنی بر ابرپیوند استفاده می‌کند.

با توجه به اطلاعات بالقوه پیکره‌های قابل مقایسه (حداقل به‌صورت عبارات زیرجمله‌ای) و اینکه کار زیادی در این زمینه برای زبان فارسی انجام نشده است روش پیشنهادی ارائه شد تا در این زمینه راه‌گشا باشد.

۳- معرفی روش پایه و پیشنهادی

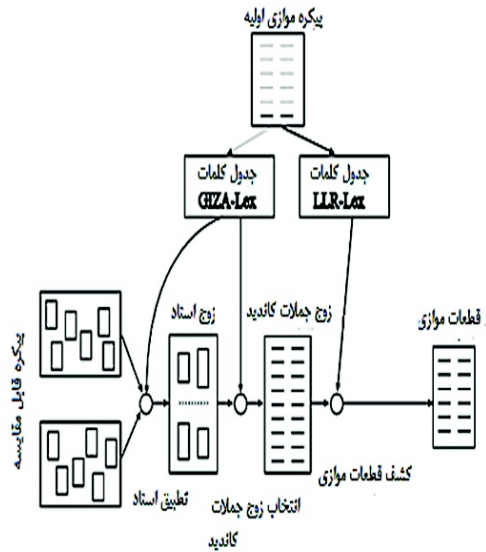
اسناد قابل مقایسه موضوع یکسانی دارند؛ اما ترجمهٔ دقیق هم نیستند. از اسناد قابل مقایسه اطلاعات مختلفی استخراج می‌شود که شامل استخراج اسناد موازی، جملات موازی و زیرجملات موازی است. به‌طور معمول پیکره‌های خیلی غیرموازی دارای زوج جملات خوب نیستند؛ یا دارای تعداد اندکی از جملات مناسب هستند و اکثر دادهٔ موازی آنها در سطح زیرجمله موجود است. در این بخش روش پایهٔ استفاده‌شده برای پروژه و روش جدید پیشنهادی معرفی می‌شود.

۳-۱- ترازبندی اسناد

برای برخی اسناد مانند اسناد استخراج‌شده از سایت khamenei.ir به‌دلیل عدم ارتباط میان اسناد فارسی و انگلیسی و نبود رابطهٔ معین، لازم است تا ابتدا ترازبندی در

۳-۲- روش پایه

به‌عنوان روش پایه از مهم‌ترین کار انجام‌شده در زمینه استخراج قطعات زیرجمله‌ای موازی یعنی روش (مانتینیو و مارکو، ۲۰۰۶) استفاده می‌شود. در این قسمت جزئیات و فرآیند پیاده‌سازی آن بیان می‌شود. فرآیند کلی برنامه در نمودار جعبه‌ای زیر مشخص است:



(شکل ۳-۲): نمودار جعبه‌ای روش پایه (مانتینیو و مارکو، ۲۰۰۶)

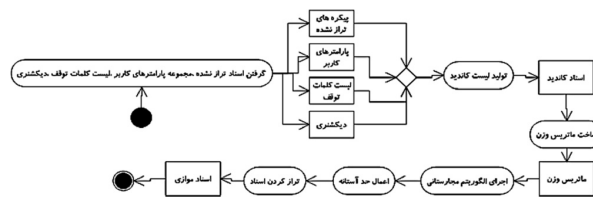
- در نخستین مرحله زوج اسناد مشابه با استفاده از ابزار بازیابی اطلاعات Lemur تشخیص داده می‌شود. به این صورت که سند مبدأ کلمه‌به‌کلمه ترجمه‌شده و به یک پرس‌وجو تبدیل می‌شود که بر روی اسناد زبان مقصد، اجرا می‌شود و بیست نتیجه بالاتر بازیابی می‌شوند. زوج‌جملات انتخاب‌شده از اسناد جفت‌شده به‌عنوان ورودی مرحله بعدی استفاده می‌شوند.
- در مرحله دوم زوج‌جملات این زوج اسناد به فیلتر انتخاب نامزد ارسال می‌شوند. در این مرحله زوج‌جمله‌هایی که تعداد کمی کلمه ترجمه‌شده دارند، حذف شده و روش کشف قطعات برای زوج جملات باقی‌مانده، اعمال می‌شود که نتیجه آن خروجی سامانه را تشکیل می‌دهد.

- در این روش از دو فرهنگ لغت دوزبانه استفاده می‌شود: جدول کلمات جیزا و جدول شباهت لگاریتمی. جدول جیزا از اجرای نرم‌افزار جیزا^۲ (پیاده‌سازی مدل‌های ترازبندی کلمات آی‌بی‌ام) روی یک پیکره موازی به‌دست می‌آید که یکی از ویژگی‌های آن این است که هر کلمه

سطح اسناد صورت گیرد. برای این منظور از روشی مبتنی بر الگوریتم مجارستانی^۱ استفاده شده و از طول و ارتباط معنایی کلمات به‌عنوان ویژگی استفاده شده است.

الگوریتم مجارستانی یک الگوریتم بهینه‌سازی ترکیبی است که مسأله تخصیص را در زمان چندجمله‌ای حل می‌کند. ورودی الگوریتم یک ماتریس غیر منفی است که درایه‌های سطر آم و ستون زام آن در حالت کلی نماینده هزینه تخصیص کار زام به فرد آم است. هدف پیدا کردن تخصیص مناسب کارها با بهینه‌ترین هزینه به افراد است. در مورد این ابزار ما با استفاده از این الگوریتم، با دادن ماتریس وزن به‌عنوان ورودی، مناسب‌ترین اسناد را به هم نگاشت می‌کنیم. پیاده‌سازی این الگوریتم جزئیات زیادی دارد، اما گام‌های آن به‌طور خلاصه برای یک ماتریس $m \times m$ به شرح زیر است:

- کوچک‌ترین عنصر هر سطر یا ستون را از تمام عناصر آن سطر یا ستون کم می‌کنیم.
 - با کمترین تعداد خط صفرهای جدول را می‌پوشانیم. اگر تعداد خطوط پوششی برابر m شد ماتریس بهینه است.
 - در غیر این صورت کوچک‌ترین عددی را که روی آن خط کشیده نشده از اعدادی که روی آنها خط کشیده نشده کم و به محل تقاطع خطوط مرحله دو اضافه می‌کنیم و دوباره صفرهای جدید را با خط پوششی، می‌پوشانیم تا به m خط برسیم.
 - مراحل ۲ و ۳ را تا رسیدن به m خط ادامه می‌دهیم. (کوهن، ۱۹۹۵)
- همچنین در برنامه از دو فایل لغت‌نامه فارسی به انگلیسی و بالعکس استفاده می‌شود.
- نمودار جعبه‌ای نحوه عملکرد این برنامه به‌صورت زیر است:



(شکل ۳-۱): نمودار فرآیند ابزار Document Aligner

² Giza++

¹ hungarian

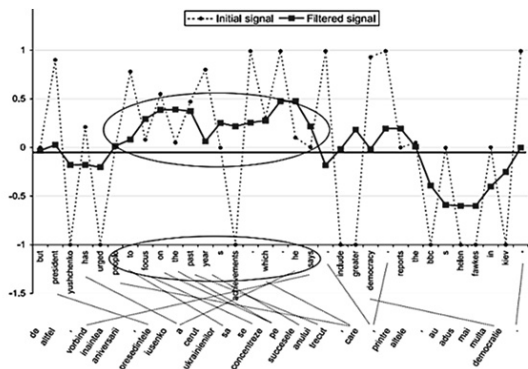
ممکن است، چندین معنا داشته باشد. جدول شباهت لگاریتمی شامل امتیازات جفت کلمات بر اساس نسبت خطی لگاریتمی^۱ یا همان LLR آنها است که با استفاده از فرمول ذکر شده در (مور، ۲۰۰۲) و روند زیر محاسبه می‌شود:

- روی پیکره موازی تراز شده، کلمات جی++ را در دو جهت اجرا می‌کنیم و هم‌ترازی‌ها را با استفاده از روش‌های مکاشفه‌ای بهبود یافته متقارن می‌کنیم.
- همه امتیازات LLR را محاسبه می‌کنیم. برای هر جفت کلمه که حداقل یک اتصال در پیکره هم‌تراز شده دارد، یک امتیاز LLR وجود دارد.
- همه $LLR(e, f)$ ها را به دسته‌های مثبت $p(e, f) > p(e).p(f)$ و منفی (ارتباط منفی) تقسیم می‌کنیم.
- برای هر f عامل نرمال‌سازی $\sum_e LLR^+(e, f)$ یا $\sum_e LLR^-(e, f)$ را محاسبه می‌کنیم.
- همه عبارات $LLR^+(e, f)$ را به عامل‌های نرمال‌سازی مربوطه تقسیم می‌کنیم تا $p^+(e|f)$ را به دست آوریم.
- همه عبارات $LLR^-(e, f)$ را به عامل‌های نرمال‌سازی مربوطه تقسیم می‌کنیم تا $p^-(e|f)$ را به دست آوریم.

دلیل استفاده از مفهوم دسته‌های مثبت و منفی این است که با توجه به فرهنگ لغت احتمالاتی و رابطه بالا تشخیص داده شود که با احتمال زیاد یک جفت کلمه، ترجمه هم هستند یا خیر، اگر باشند مثبت و در غیر این صورت منفی دسته‌بندی می‌شوند.

پس از محاسبه جدول شباهت لگاریتمی، برای هر جفت جمله نامزد به صورت حریصانه با توجه به امتیاز LLR ترازبندی صورت گرفته و مقدار آن در نمودار قرار می‌گیرد و سیگنال تشکیل می‌شود. قرارگیری مقادیر در نمودار به این صورت است که برای هر جفت کلمه اگر p^+ وجود داشته باشد، یعنی تخصیص مثبتی بین جفت‌های حاصل از کلمه مبدأ و هر یک از کلمات جمله مقصد وجود خواهد داشت و بالاترین مقدار p^+ ها در نمودار قرار می‌گیرد و اگر p^+ وجود نداشته باشد و همه تخصیص‌ها منفی باشد، کمترین مقدار منفی‌ها در نمودار قرار می‌گیرد. این امر به این دلیل است که کلمه‌ای که هیچ یک از تخصیص‌های مثبت نیستند، به حتم ترجمه بهتری دارد که در این جمله نیست و نباید ترازبندی برایش صورت بگیرد. پس کم‌ترین عدد در بین p^- ها برایش انتخاب می‌شود تا در نرمال‌سازی حذف شود. اگر

هم کلمه‌ای در جمله مقصد هیچ گونه ترجمه‌ای، حتی با p^- در جمله مقصد نداشت، مقدار آن در نمودار ۱- قرار می‌گیرد. پس از این مرحله بین مقادیر همسایه، میانگین‌گیری صورت می‌گیرد (هر مقدار با دو مقدار مجاورش) و مقادیر نرمال‌سازی می‌شوند. پس از نرمال‌سازی، بخش‌هایی از نمودار که مقادیر متوالی مثبت دارند، به عنوان قطعه استخراج می‌شوند. شکل (۳-۳) نمودار سیگنال قبل و بعد از فیلتر شدن و قطعه استخراج شده را نشان می‌دهد.



(شکل ۳-۳): نمودار سیگنال برای یک جفت جمله (مانتینیو و مارکو، ۲۰۰۶)

۳-۳-۳- روش پیشنهادی

۳-۳-۳-۱- توضیح کلی الگوریتم

روش پیشنهادی، روشی برای استخراج قطعات موازی از پیکره قابل‌مقایسه است. در این روش نیز همانند روش پایه از دو فرهنگ لغت استفاده می‌شود که نخستین‌شان فرهنگ لغت ساده برای تعیین جملات نامزد و دیگری فرهنگ لغت مبتنی بر LLR است؛ که در قسمت تعیین قطعات استفاده می‌شود. فرهنگ لغت مبتنی بر LLR با استفاده از روش ذکر شده در روش پایه ایجاد می‌شود.

در مرحله نخست بر روی متون ورودی پیش‌پردازش انجام می‌شود. در بخش پیش‌پردازش، زیربخش‌های توکن‌بندی، نرمال‌سازی و تعیین مرز جملات صورت می‌گیرد. بدین شکل که نویسه‌های متن ورودی یکسان‌سازی، فواصل و خطوط اضافی حذف و در مرحله بعدی به واحدهای جمله شکسته می‌شوند. در این بخش مرز جمله‌ها با استفاده از علائم جداکننده جمله شامل «»، «»، «؟» و «؟» و مرز کلمه‌ها با استفاده از علائم فضای خالی،

^۱ Log Linear Ratio

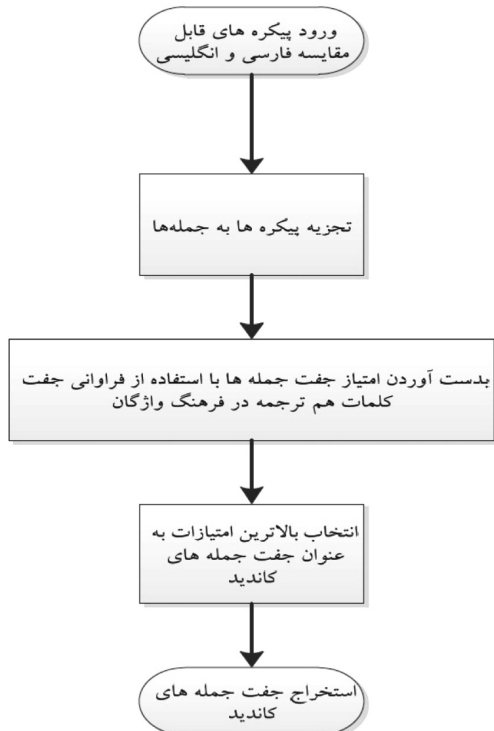
می‌شود). این کار به دلیل حفظ ساختار قطعه صورت می‌گیرد.

➤ امتیاز به‌روزرسانی می‌شود.

اگر امتیاز بلوک از حد خاصی افت کند، این جستجو خاتمه پیدا می‌کند. همچنین هر قطعه اگر انتخاب شود، کلمات آن در قطعات بعدی لحاظ نخواهند شد (قطعات هم‌پوشانی ندارند). شکل (۳-۵) مثالی از روند اعمال الگوریتم به یک جمله و گسترش بلوک مسیر حول یک نقطه پررنگ را نشان می‌دهد.

شکل (۳-۶) نمودار جعبه‌ای الگوریتم را نشان می‌دهد.

در شکل (۳-۵) مثالی از نحوه عملکرد الگوریتم پیشنهادی آورده شده است. در این شکل یک جفت جمله نامزد از پیکره قابل مقایسه انتخاب شده است که روی دو محور قرار گرفته‌اند. همان‌طور که در شکل مشخص است، جفت کلمات (یک و one) در مرحله نخست بالاترین امتیاز LLR را داشته‌اند و بلوک استخراج قطعه حول آن نقطه تشکیل شده است. در هر مرحله مطابق الگوریتم با توجه به امتیازات نقاط همسایه بعدی انتخاب شده و مسیر در بلوک گسترده می‌شود تا در نهایت با گذشتن از حد آستانه جستجو خاتمه یابد.



(شکل ۳-۴): نمودار جعبه‌ای به‌دست آوردن جملات نامزد

«»، «>»، «<»، «[»»، «]»»، «-»، «>»، «<» و «/» و خط جدید مشخص می‌شود.

پس از این مرحله مجموعه جملات نامزد مشخص می‌شوند. بدین صورت که امتیاز تمامی جفت جمله‌ها با استفاده از فراوانی جفت کلمات هم ترجمه موجود در فرهنگ واژگان به‌دست می‌آید و بالاترین امتیازات به‌عنوان جفت جمله‌های نامزد انتخاب می‌شوند. شکل (۳-۴) فرآیند بخش انتخاب جملات نامزد را نشان می‌دهد.

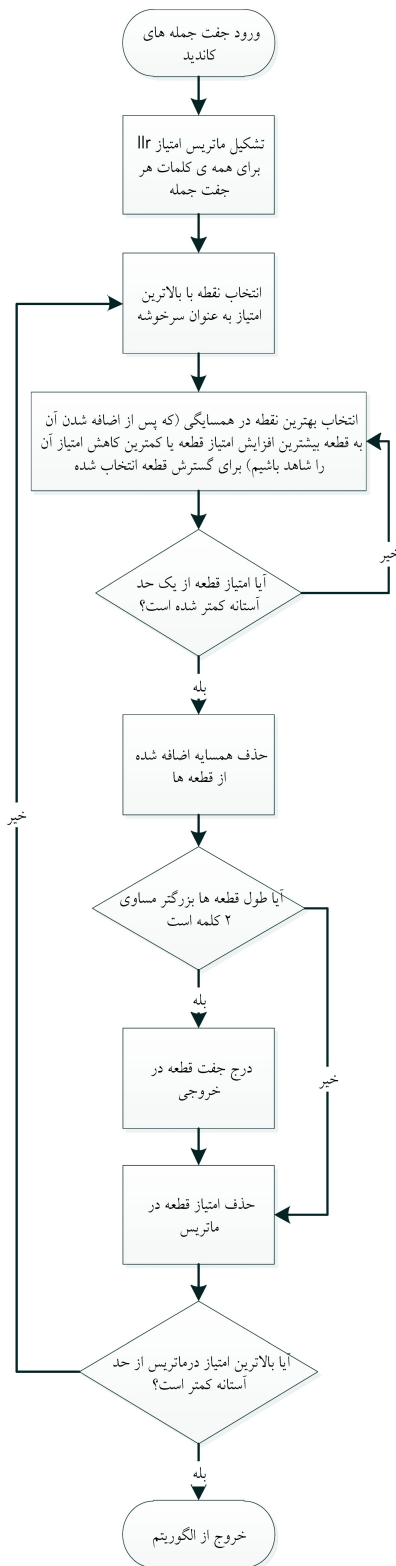
در روش پیشنهادی برای در نظر گرفتن حالت‌های مختلف ترازبندی، از ماتریس امتیازات استفاده می‌کنیم. به این صورت که برای هر جفت جمله نامزد، ماتریس امتیازات LLR متناظر جفت کلمات ایجاد می‌گردد؛ سپس به‌صورت نزولی از نقطه دارای بیش‌ترین امتیاز، به‌عنوان نقطه شروع استخراج قطعه آغاز شده و پررنگ‌ترین نقاط یعنی نقاط دارای بالاترین امتیاز تا وقتی که از حد آستانه‌ای افت نکرده‌اند، برای استخراج قطعه، یکی پس از دیگری پردازش می‌شوند. این حد آستانه به‌صورت تجربی به‌دست می‌آید؛ چون با پایین آوردن آن کیفیت قطعه‌ها پایین می‌آید؛ پس تا زمانی که کیفیت قطعات مطلوب باشد، حد آستانه را پایین می‌آوریم. برای هر نقطه پررنگ در واقع خوشه‌بندی صورت می‌گیرد و هر خوشه نمایان‌گر یک قطعه است.

فرآیند پیدا کردن قطعه برای هر نقطه پررنگ به شرح زیر است:

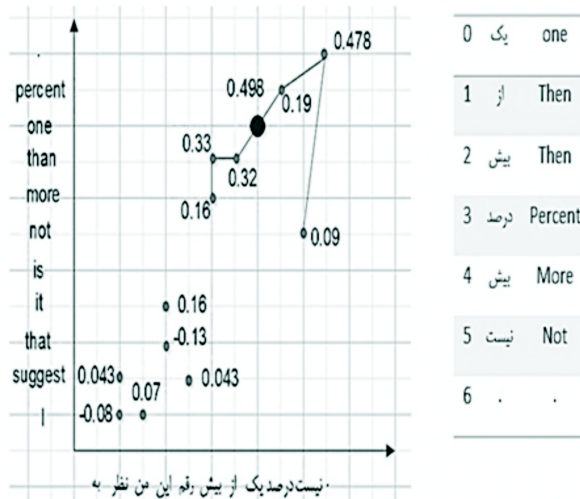
➤ امتیاز قطعه با توجه به ویژگی‌ها و الگوها محاسبه می‌شود. الگوها با توجه به نوع زبان‌های استفاده‌کننده از برنامه (فعالاً فارسی-انگلیسی) و ساختار معمول قطعات در این نوع زبان‌ها که از پیکره آموزشی یاد گرفته شده‌اند، انتخاب شدند. به‌عنوان مثال ساختار قطری با شیب ۱ یا ۱- یکی از ساختارهای بسیار رایج قطعات در زبان فارسی و انگلیسی است. بدین معنی که کلمات یا به‌ترتیب (شیب ۱) و یا به‌طور معکوس (شیب ۱-) با هم متناظرند. ویژگی‌های مورد استفاده نیز در ادامه توضیح داده خواهند شد.

➤ در همسایگی قطعه (بار نخست قطعه تنها یک نقطه است و همسایگی آن هشت نقطه مجاور آن است) جستجو می‌کنیم و امتیاز کل قطعه را در صورت افزایش هر کلمه همسایه محاسبه کرده و کلمه‌ای که منجر به بالاترین امتیاز شده انتخاب و به مسیر افزوده می‌شود.

➤ کلمه انتخاب‌شده به کلمه‌ای که x یا y آن کم‌ترین فاصله را با آن دارد، وصل می‌شود (والدش انتخاب



شکل ۳-۶: روندنمای الگوریتم پیشنهادی



بیش از یک درصد نیست. ⇔ not more than one percent.

شکل ۳-۵: مثالی از شیوه اعمال الگوریتم پیشنهادی

۳-۳-۲- ویژگی‌های استفاده شده در روش پیشنهادی

۱- طول قطعه (F1): هر چه طول قطعه بیشتر باشد بهتر است و امتیازها به نسبت طول مقایسه می‌شوند.

$$F1 = \text{Path Length} \quad (1)$$

تعداد تخصیص‌ها = Path Length

۲- میانگین LLRها (F2): LLR عامل نخستین است، به هم ترجمه بودن کلمات قطعه اشاره دارد؛ پس مقدار آن هر چه بیشتر باشد، بهتر است.

$$F2 = \frac{\sum_{i=1}^N LLR_i}{N} \quad (2)$$

۳- طول مسیر در قطعه (F3): بین کم‌ترین و بیش‌ترین اندیس انتخاب‌شده از ماتریس، هر چه تعداد کلمات دارای معادل بیشتر باشد، بهتر است و امتیاز بالاتری به آن تخصیص می‌یابد.

$$F3 = \frac{I}{X} \quad (3)$$

ا: تعداد تخصیص‌ها

X: حداکثر طول بین دو جمله فارسی و انگلیسی

۳- میزان شیب اسکلت قطعه در بلوک (F4): گفتیم که برای هر قطعه اسکلت مسیر تشکیل می‌شود. هر چه قدرمطلق شیب این مسیر به یک نزدیک‌تر باشد، مسیر به بهینه‌ترین الگوی قطعات فارسی-انگلیسی یعنی همان حالت قطری نزدیک‌تر است. در قطعات هم‌ترجمه انگلیسی-فارسی بیش‌ترین الگوی رخ داده، حالت قطری با شیب ۱ یا ۱- است. (به‌عنوان مثال ترازبندی کلمات دو عبارت "دوست خوب" و "friend"good در نمودار کلمات آنها که محور افقی عبارت انگلیسی و محور عمودی عبارت فارسی را نشان بدهد، به‌صورت یک نمودار خطی با شیب ۱- است). برای محاسبه این ویژگی، قدرمطلق شیب خطوط بین هر دو نقطه در مسیر قطعه محاسبه می‌شود.

- اگر کوچک‌تر از ۱ است: خودش

- اگر بزرگ‌تر از ۱ است: معکوسش

سپس از مقادیر حاصله میانگین می‌گیریم. هر چه این مقدار به یک نزدیک‌تر باشد، بهتر است. (مقادیر افقی و عمودی نیز حذف می‌شوند)

۴- پراکندگی شیب در خطوط بین نقاط تشکیل‌دهنده مسیر در بلوک قطعه (F5): شکستگی مسیر را بررسی می‌کند. هر چه شکستگی کمتر باشد، مسیر به حالت قطری نزدیک‌تر است (شکستگی بیشتر از یک به‌طور معمول به قطعات مناسبی منجر نمی‌شود).

- جدول توزیع شیب‌ها به دست می‌آید (شیب‌ها و تعداد تکرارشان).

- میزان شیب، عبارت از تعداد درایه‌های جدول تقسیم بر تعداد یال‌های مسیر است.

۵- مربعی بودن بلوک قطعه (F6): هر چه بلوک مربعی‌تر باشد، یعنی طول قطعه فارسی و انگلیسی معادل به هم نزدیک‌تر است.

$$F6 = \left| \frac{\max x - \max y}{\max y - \min y} \right| \quad (4)$$

Maxx: بیشینه اندیس‌ها در قطعه فارسی

Minx: کمینه اندیس‌ها در قطعه فارسی

Maxy و Miny نیز به همین ترتیب برای قطعه انگلیسی)

اگر مقدار A بزرگ‌تر از یک بود، خودش و اگر کمتر از یک بود معکوسش به‌عنوان امتیاز مربعی بودن بلوک در نظر گرفته می‌شود.

این امتیازات طوری تراز شده‌اند که همگی بین ۰ و ۱۰ قرار گیرند و با استفاده از ترکیب خطی به‌صورت وزن‌دار ترکیب شده‌اند و امتیاز نهایی بلوک را شکل می‌دهند.

$$Score = \frac{\sum_{i=1}^6 w_i \cdot f_i}{\sum_{i=1}^6 w_i} \quad (5)$$

۳-۳-۳- وزن‌دهی ویژگی‌ها

برای تعیین وزن ویژگی‌ها، یک پیکره طلایی شامل هشتصد جمله از پیکره‌های Tehran avenue، central asia و wikipedia و ۱۶۳۰ قطعه مشخص شده در آنها به‌عنوان مرجع ساخته شد و با استفاده از روش تبرید شبیه‌سازی شده^۱ (یا تبرید تدریجی فلزات) بهینه‌ترین وزن‌ها تخمین زده شد.

تبرید شبیه‌سازی شده، یک روش بهینه‌سازی فراابتکاری ساده و مفید در حل مسائل بهینه‌سازی است. تکنیک تبرید تدریجی، به‌وسیله متالورژیست‌ها برای رسیدن به حالتی که در آن ماده جامد، به‌خوبی مرتب و انرژی آن کمینه شده باشد، استفاده می‌شود. این تکنیک شامل قراردادن ماده در دمای بالا و سپس کم کردن تدریجی این دماست. (کرکپاتریک و همکاران، ۱۹۸۳)

در روش شبیه‌سازی تبریدی، هر نقطه در فضای جستجو مشابه یک حالت از یک سامانه فیزیکی است؛ و انرژی داخلی سامانه در آن حالت باید کمینه شود. در این روش، هدف انتقال سامانه از حالت نخستین دلخواه، به حالتی است که سامانه در آن کمترین انرژی را داشته باشد.

برای حل یک مسئله بهینه‌سازی، الگوریتم ابتدا از یک جواب نخستین شروع و سپس در یک حلقه تکرار به جواب‌های همسایه حرکت می‌کند. اگر جواب همسایه بهتر از جواب فعلی باشد، الگوریتم آن را به‌عنوان جواب فعلی قرار می‌دهد (به آن حرکت می‌کند)، در غیر این صورت، الگوریتم آن جواب را با احتمال $\exp(-\Delta E/T)$ به‌عنوان جواب فعلی می‌پذیرد. در این رابطه ΔE تفاوت بین تابع هدف جواب فعلی و جواب همسایه است و T یک پارامتر به نام دماست. در هر دما، چندین تکرار اجرا و سپس دما به‌آرامی کاهش داده می‌شود. در گام‌های نخستین دما خیلی بالا قرار داده می‌شود تا احتمال بیشتری برای پذیرش جواب‌های بدتر وجود داشته باشد. با کاهش تدریجی دما، در گام‌های پایانی احتمال کم‌تری برای پذیرش جواب‌های بدتر وجود خواهد

¹ Simulated annealing

داشت؛ و بنابراین الگوریتم به سمت یک جواب خوب، همگرا می‌شود.

در هر مرحله، الگوریتم تبرید شبیه‌سازی شده، چند حالت را در همسایگی حالت کنونی در نظر می‌گیرد و به‌طور احتمالی تصمیم می‌گیرد که سامانه را از حالت منتقل کند یا در همین حالت باقی بماند. این احتمالات در نهایت سامانه را به حالت با انرژی کمتر میل می‌دهد. (کرنی، ۱۹۸۵) نمودار جعبه‌ای این الگوریتم در شکل (۷-۳) آورده شده است.

اعمال الگوریتم تبرید شبیه‌سازی شده

برای تعیین بهینه وزن‌های ویژگی‌ها، آنها را به‌صورت یک بردار شش تایی در نظر گرفتیم که هر بار به‌صورت تصادفی شش مقدار آن تغییر می‌کنند. به‌صورتی که به احتمال ۳۳٪ خودش باقی می‌ماند، به احتمال ۳۳٪/۰/۱ افزایش می‌یابد و یا به احتمال ۳۳٪ به میزان ۰/۱ کاهش می‌یابد. به این صورت به نقطه همسایه حرکت می‌کند.

در این مرحله برنامه فراخوانی می‌شود و دقت و بازخوانی با مقایسه با پیکره طلایی مرجع محاسبه می‌شود. هزینه نیز در پیاده‌سازی انجام‌شده به‌عنوان معکوس امتیاز در نظر گرفته می‌شود و ΔE برابر با هزینه نقطه جدید منهای هزینه نقطه قبلی است.

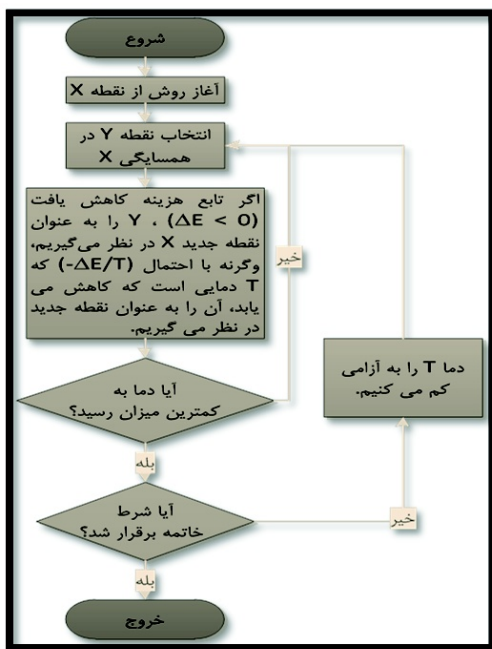
اگر این مقدار کمتر از صفر بود، همسایه به‌عنوان نقطه جدید در نظر گرفته می‌شود و اگر نبود به احتمال P که در الگوریتم ذکر شد، نقطه همسایه به‌عنوان نقطه جدید قبول می‌شود. در هر بار اجرا دما یک درجه کاهش می‌یابد و زمانی که دما صفر شود، جواب برگردانده می‌شود.

۳-۳-۴- نمونه‌ای از خروجی سامانه

در جدول زیر تعدادی از عبارات استخراج‌شده توسط سامانه پیشنهادی آورده شده است.

۴- ارزیابی

هدف از این بخش ارزیابی سامانه پیشنهادی در مقایسه با سامانه پایه با استفاده از معیارهای ارزیابی است. همچنین بررسی این موضوع که پیکره موازی استخراج‌شده توسط ابزار پیشنهادی چه تأثیری بر کیفیت سامانه ترجمه ماشینی دارد.



(شکل ۳-۷): دیاگرام بلوکی الگوریتم تبرید شبیه‌سازی شده (کرنی، ۱۹۸۵)

(جدول ۳-۱): نمونه‌ای از قطعات استخراج‌شده توسط سامانه پیشنهادی

announced that the government will	-> اعلام کرد که دولت
is that many of these	-> این است که برخی از این
concert that night was supposed	-> که قرار بود همان شب کنسرت
of love and fear and	-> از عشق و ترس و
foreign trade	-> تجارت خارجی
believes the	-> معتقد است
the first step in	-> اولین گام در
presidential election	-> انتخابات ریاست جمهوری
these works are either	-> این آثار یا
limitations in the use	-> محدودیت در استفاده
young soldier	-> سرباز جوان
is that the index	-> این است که فهرست
the country's economic development	-> توسعه اقتصادی این کشور
journalist and analyst	-> روزنامه نگار و تحلیلگر
corrupt security forces	-> نیروهای امنیتی فاسد
really effective	-> واقعا مؤثر

(جدول ۴-۱): مشخصات پیکره موازی و فرهنگ لغت‌های

استفاده‌شده

پیکره موازی اولیه	PEN دارای ۳۰۴۷۹ جفت جمله موازی
فرهنگ لغت ساده	دارای ۷۱۰۹۰ درایه
(جدول کلمات LLR)	دارای ۲۲۹۷۴۷ درایه

همچنین وزن‌های تخمین زده شده در جدول زیر آورده شده‌اند:

(جدول ۴-۲): وزن‌های تخمین زده شده برای ویژگی‌ها

ویژگی‌ها	وزن‌های تخمینی
طول قطعه	۰/۱۳
میانگین LLR	۰/۱۵
پراکندگی شیب مسیر در بلوک قطعه	۰/۴۵
درصد وجود کلمات هم ترجمه در قطعه	۰/۹۲
مربعی بودن بلوک قطعه	۰/۵۱
میزان شیب مسیر در بلوک قطعه	۰/۴۴

برای آزمون سامانه حد آستانه مربوط به بخش امتیاز سامانه که با افت کردن امتیاز از آن، جستجو برای قطعه پایان می‌یابد، با توجه به محدوده امتیازات (۰-۱) و وزن‌های به دست آمده، به صورت تجربی با مشاهده کیفیت قطعات حاصل، پنج قرار داده شد.

محاسبه دقت، بازخوانی و معیار F بر مبنای کلمات صورت گرفت. به این صورت که دقت هر قطعه عبارت از تعداد کلمات استخراج شده صحیح نسبت به تعداد کلمات استخراج شده و بازخوانی عبارت از تعداد کلمات استخراج شده صحیح نسبت به تعداد کلماتی که مطابق مرجع باید استخراج می‌شده، است. به عنوان مثال:

مرجع: گل سرخ تازه => fresh red flower
استخراج شده توسط سامانه: گل سرخ => red flower

دقت = ۱۰۰٪

بازخوانی = ۶۶٪

۴-۱- معیارهای ارزیابی

مشهورترین معیارهای ارزیابی در حوزه ترجمه ماشینی عبارتند از دقت، بازخوانی، معیار F و بلو که ما از این معیارها در ارزیابی سامانه استفاده خواهیم کرد.

۴-۱-۱- معیار F^۱

در حوزه پردازش زبان طبیعی، عبارت معیار F به ترکیبی از دقت و بازخوانی^۲ اشاره دارد. این معیار به طور عمومی برای ارزیابی سامانه‌های استخراج اطلاعات استفاده می‌شوند. فرض کنید مجموعه نامزد با Y و مجموعه مراجع را با X نمایش دهیم. در این صورت سه مقدار دقت، بازخوانی و معیار F به صورت زیر تعریف می‌شوند:

$$\text{recall}(Y|X) = \frac{|X \cap Y|}{X} \quad (6)$$

$$\text{precision}(Y|X) = \frac{|X \cap Y|}{Y} \quad (7)$$

$$F\text{-measure} = \frac{2 * \text{precision} * \text{recall}}{\text{precision} + \text{recall}} \quad (8)$$

۴-۱-۲- معیار بلو^۳

معیار بلو برای ارزیابی خودکار ترجمه ماشینی از مقایسه خروجی سامانه با ترجمه‌های صحیح استفاده می‌کند (پایینی و روکوس، ۲۰۰۱). بنابراین، به یک پیکره آزمایشی که حداقل شامل یک ترجمه دستی از هر جمله است نیاز دارد. در زمان آزمایش، هر جمله آزمایشی به سامانه ترجمه ماشینی داده می‌شود و خروجی با مقایسه با ترجمه درست امتیازدهی می‌شود. این امتیاز را امتیاز بلو می‌نامند. جمله خروجی را جمله نامزد^۴ و ترجمه‌های درست را مراجع^۵ می‌نامند. امتیاز بلو به کمک دو عامل دقت^۶ و طول جمله نامزد محاسبه می‌شود.

۴-۲- شرایط آزمون

۴-۲-۱- پارامترهای سامانه پیشنهادی

برای آموزش نخستین سامانه و ساخت (جدول کلمات LLR، از فرهنگ لغت احتمالاتی مستخرج از پیکره PEN استفاده شده است که یک پیکره موازی خبری است.

- 1 F-measure
- 2 Recall
- 3 BLEU
- 4 Candidate Sentence
- 5 References
- 6 Precision

۴-۲-۲- پیکره‌های آزمون

برای آزمایش سامانه یک پیکره طلایی آزمون شامل ۱۸۰۰ زوج جمله و قطعات مرجع معین در آن با شرایط ذکر شده در جدول (۳-۴) به سامانه داده شد. لازم به ذکر است، جملات انتخاب شده برای بخش آزمون از دو پیکره Tehran avenue و central asia با جملات موجود در پیکره تنظیم وزن‌های سامانه، هم‌پوشانی نداشتند.

(جدول ۴-۴): تعداد اسناد تراز شده

تعداد اسناد تراز شده	ویژگی‌های استفاده شده
۶۷۲	طول
۶۱۲	ترجمه
۵۸۲	طول+ترجمه
۴۷۶	طول+ترجمه+تاریخ

(جدول ۴-۳): مشخصات دادگان ارزیابی

نام پیکره	تعداد زوج جملات	موضوع پیکره	کیفیت ترجمه پیکره	تعداد قطعات مرجع
Khameneii.ir	۵۰۰ زوج جمله	سیاسی	ضعیف- متون مشترک کم	۴۷۵
Tehran avenue	۵۰۰ زوج جمله	هنری- عمومی	متوسط- قابل مقایسه	۶۰۲
Central asia	۵۰۰ زوج جمله	سیاسی- خبری	خوب-نیمه موازی	۶۵۲
Books	۳۰۰ زوج جمله	عمومی	خوب-قابل مقایسه	۲۷۸

همان‌طور که نتایج نشان می‌دهد، استفاده از تنها یک ویژگی طول یا ترجمه برای ترازبندی اسناد کافی نیست. استفاده تنها از مدل طول، نخستین سندی را که از نظر طول با سند مبدأ با توجه به نسبت طول تطابق داشته باشد، را به آن تخصیص می‌دهد و استفاده تنها از مدل ترجمه، ممکن است اسنادی را با اختلاف طول زیاد فقط به دلیل داشتن عبارات مشترک به هم تخصیص دهد. به این دلیل تعداد اسناد خروجی در آزمایش استفاده از این دو ویژگی به تنهایی بالا، اما دقت آنها پایین است. استفاده توأم از این دو ویژگی منجر به نتایج منطقی می‌شود؛ اما چون مجموعه نامزد بزرگ است، زمان رویه بسیار طولانی خواهد شد. استفاده از ویژگی تاریخ، مجموعه نامزد را هرس می‌کند و در نتیجه در زمان کوتاهی نتایج دقیقی ارائه می‌کند که در جدول (۴-۵) قابل مشاهده است. واضح است که ویژگی تاریخ تنها می‌تواند به‌عنوان یک ویژگی جانبی استفاده شود و استفاده از آن بدون دو ویژگی اصلی طول و ترجمه بی‌معنی است.

برای محاسبه دقت، بازخوانی و معیار F، ۱۰۰ تا از جفت اسناد خروجی به‌صورت کاملاً تصادفی انتخاب و صحت تخصیص آنها بررسی شد. در مورد بازخوانی پس از انجام بررسی‌ها مشخص شد ۹۱ مورد از جفت‌اسناد به‌طور واقعی موازی بوده‌اند.

۴-۳- روش ارزیابی

برای ارزیابی، سامانه را با توجه به نوع و کیفیت پیکره ورودی، با ۳ معیار دقت، بازخوانی و معیار F در دو سامانه پایه و سامانه پیشنهادی بررسی می‌کنیم. همچنین تأثیر دادگان استخراج شده توسط سامانه پیشنهادی را در سامانه ترجمه ماشینی پایه با استفاده از معیار بلو می‌سنجیم.

۴-۴- مقایسه و ارزیابی سامانه پایه و سامانه

پیشنهادی

۴-۴-۱- بررسی عملکرد سامانه پیشنهادی

• ارزیابی بخش ترازبندی اسناد

برای ارزیابی این برنامه ۶۷۲ سند انگلیسی و ۵۵۰۰ سند فارسی از مجموعه اسناد سایت khameneii.ir در نظر گرفته شد و اسناد با توجه به مدل ترجمه، مدل طول، هر دو ویژگی و همچنین هر دو ویژگی به‌علاوه پیشوند تاریخ ترازبندی شدند. جدول (۴-۴) تعداد زوج‌اسناد تراز شده خروجی از اسناد ورودی در هر آزمایش را نشان می‌دهد.

(جدول ۴-۵): ارزیابی Document Aligner

معیار ارزیابی	معیار F	بازخوانی	دقت
طول	۴۴	۴۶	۴۲
ترجمه	۶۶	۷۰	۶۴
طول+ترجمه	۷۹	۸۳	۷۶
طول+ترجمه+تاریخ	۸۸	۹۳	۸۴

• ارزیابی بخش اصلی سامانه پیشنهادی

در این بخش عملکرد سامانه پیشنهادی ارزیابی می‌شود. این ارزیابی در دو بخش صورت می‌گیرد:

- ۱- پیکره استخراج شده توسط برنامه چقدر خوب است.
- ۲- روی سامانه ترجمه ماشینی تا چه اندازه تأثیرگذار است.

در این راستا ابتدا با استفاده از سه معیار دقت، بازخوانی و معیار F برای هر یک از چهار گروه داده آزمون سنجیده شده است؛ سپس سامانه پیشنهادی روی سامانه ترجمه ماشینی پایه آزمایش شده است.

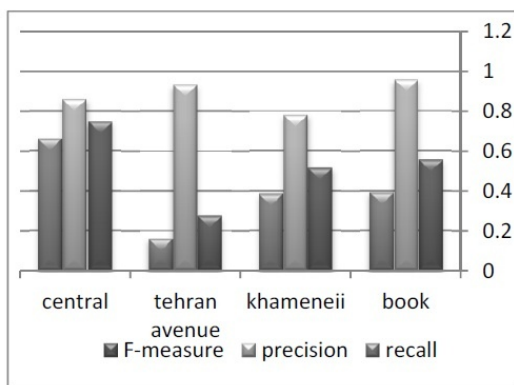
پیکره khameneii.ir یک پیکره نوفه‌ای با کیفیت ترجمه نامطلوب است که حوزه آن سیاسی محسوب می‌شود. برای این پیکره به دلیل سخنرانی بودن اغلب متون و ثقیل بودن آنها، اغلب کلمات در جدول کلمات LLR یافت نشده و بازخوانی سامانه پایین است؛ اما قطعه‌هایی که استخراج شده‌اند کیفیت تاحدودی خوب دارند. به دلیل محاوره‌ای بودن متون، قطعات زیاد از ساختار معمول زبان پیروی نکرده و باعث کاهش دقت الگوریتم نسبت به سایر پیکره‌ها شده است.

پیکره Tehran Avenue یک پیکره به‌طور کامل قابل‌مقایسه در حوزه عمومی است که کیفیت ترجمه متوسطی داشته و دارای اسامی خاص بسیاری است که به‌طور طبیعی در جدول LLR موجود نیستند؛ در نتیجه بازخوانی مربوط به این پیکره بسیار پایین است. از سوی دیگر کلماتی که در جدول موجود بوده‌اند با دقت زیادی قطعه تولید کرده‌اند که این نشان‌دهنده دقت الگوریتم است. پیکره central asia نیمه‌موازی، با کیفیت ترجمه بالا و در حوزه خبری است. به دلیل اینکه این حوزه با حوزه فرهنگ لغت استفاده شده برای ساخت جدول LLR به‌طور کامل مطابق است، مشاهده می‌شود که بازخوانی الگوریتم برای این پیکره، بالا و دقت آن بسیار بالا و در کل کارایی سامانه مطلوب است.

پیکره Book شامل متون عمومی و اسامی خاص زیاد است و کیفیت ترجمه متوسط است. در جملات ورودی این پیکره جملات کوتاه و محاوره‌ای نیز مشاهده می‌شود. به دلیل عدم اشتراک حوزه این پیکره با پیکره استفاده شده برای ساخت جدول LLR بازخوانی برنامه برای این پیکره نیز تاحدودی پایین است. همچنین کوتاه بودن جملات به دلیل عدم امکان مانور الگوریتم برای استخراج قطعه‌ها مزید بر

علت است، اما دقت برنامه به دلیل ترجمه تاحدودی خوب کتاب‌ها برای لغات موجود در جدول LLR بالاست.

در شکل (۴-۱) عملکرد الگوریتم پیشنهادی به صورت کلی و به صورت مقایسه‌ای برای پیکره‌های مختلف آورده شده است. با توجه به نمودار، واضح است، مهم‌ترین نکته برای افزایش کارایی برنامه، تطابق حوزه پیکره آزمون با حوزه متون آموزشی برای ساخت جدول LLR است، که باعث می‌شود تعداد عبارات استخراج شده و در نتیجه بازخوانی برنامه بالا برود. همچنین ساختار جملات ورودی در دقت برنامه مؤثر است. چون الگوهای قطری که به‌عنوان ویژگی در برنامه استفاده شده‌اند، از روی ساختار زبانی برداشته شده‌اند. مشاهده می‌شود که دقت کلی برنامه بسیار بالاست؛ که این امر به دلیل دقت در انتخاب ویژگی‌های مناسب است. اما نکته‌ای که در این نمودار به چشم می‌آید، بازخوانی پایین نرم‌افزار است. با بررسی دادگان خروجی مشاهده شد تعداد زیادی از قطعه‌های استخراجی و مرجع برای پیکره‌های Tehran avenue و book طول دو دارند. این شبهه ایجاد شد که شاید الگوریتم از نظر ساختاری به داده‌های با طول کم متمایل شده باشد. در این راستا قطعه‌های با طول ۲ از مرجع حذف و بهینه‌سازی پارامترها دوباره انجام شد. این بار وزن ویژگی طول حدود ۰/۲ افزایش یافت. دقت و بازخوانی دوباره محاسبه شد.



(شکل ۴-۱): ارزیابی سامانه پیشنهادی روی کل پیکره‌ها

مشاهده می‌شود که این بار بازخوانی برنامه بالاتر رفته است؛ اما از دقت الگوریتم کاسته شده است. این امر به این دلیل است که قطعه‌ای استخراجی به سمت طول بیشتر هدایت شده‌اند که این کار با افزایش تعداد کلمات استخراجی، احتمال رخداد کلمات اشتباه را نیز بالا می‌برد.

(جدول ۴-۶): نتایج ارزیابی الگوریتم پیشنهادی روی سامانه

ترجمه		
افزوده	سامانه پایه	سامانه با داده‌های
۱۸.۴۶	۱۸.۱۶	مجموعه آزمون ۱
۱۵.۷۵	۱۴.۳۶	مجموعه آزمون ۲

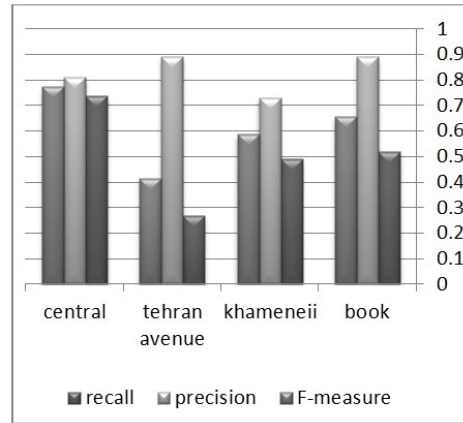
همان‌طور که مشاهده می‌شود، افزودن دادگان استخراج‌شده توسط الگوریتم، اگرچه دارای حجم ناچیزی به نسبت داده‌های پیکره PEN بودند (نسبت ۱ به ۱۵) در هر دو حالت آزمون در بلوی سامانه بهبود ایجاد کرده است. در حالت اول که مجموعه آزمون از خود پیکره PEN انتخاب شده بود این بهبود کمتر و برابر ۰/۳۳ است؛ اما در حالت دوم که بخشی از مجموعه آزمون از khameneii.ir انتخاب شده بود اگر چه بلوی سامانه کاهش پیدا کرد، اما به‌صورت مقایسه‌ای نسبت به سامانه پایه ۱/۴ واحد بلو افزایش داشته است. این امر به این دلیل است که درون دادگان آموزشی سامانه دو، دادگان استخراج‌شده از khameneii.ir وجود داشته است که به ترجمه بهتر مجموعه آزمون کمک کرده است. دلیل کاهش بلوی کلی سامانه‌ها نیز این است که بیشتر داده آموزشی از پیکره PEN بوده است و حجم داده آموزشی مناسب برای مجموعه آزمون کافی نبوده است. همچنین متون آزمون ساختار متن‌ی سخنرانی و سخت داشته که ترجمه آن برای سامانه ساده نبوده است.

در (جدول ۴-۷) نمونه‌هایی از جملات ترجمه‌شده توسط دو سامانه ترجمه ماشینی مورد ارزیابی (سامانه ترجمه پایه و سامانه ترجمه با دادگان افزوده استخراج‌شده) برای تأیید صحت بهبود آورده شده است.

۴-۴-۳- مقایسه سامانه‌های پایه و پیشنهادی

در این بخش به مقایسه کارایی سامانه پایه، سامانه پایه بهبودیافته و سامانه پیشنهادی با استفاده از سه معیار دقت، بازخوانی و معیار F بر روی پیکره‌های آزمون پرداخته می‌شود.

همان‌طور که در شکل (۴-۳) مشاهده می‌شود، دقت روش پیشنهادی به نسبت روش پایه، تفاوت فاحشی دارد. این امر به دلیل استفاده از ویژگی‌های عنوان‌شده و در نظر گرفتن الگوهای مناسب زبانی برای استخراج بلوک‌های قطعه است.



(شکل ۴-۲): ارزیابی مجدد سامانه پیشنهادی روی کل پیکره‌ها

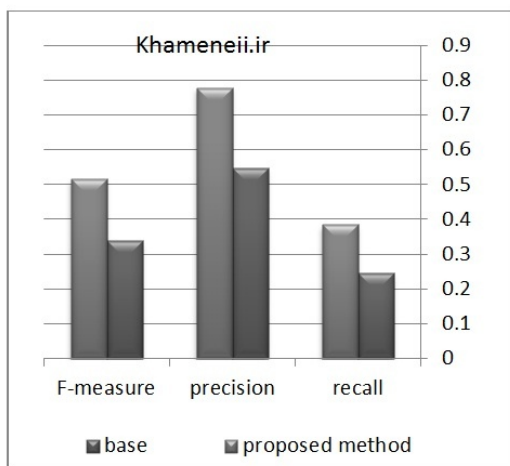
۴-۴-۲- بررسی تأثیر دادگان استخراج‌شده بر روی سامانه ترجمه ماشینی

در بخش دوم، ارزیابی تأثیر دادگان استخراج‌شده را بر روی سامانه ترجمه ماشینی بررسی می‌کنیم. الگوریتم پیشنهادی بر روی پیکره‌های قابل‌مقایسه Tehran avenue، khameneii.ir و book اعمال شد که پیش از این با استفاده از روش‌های موجود به دلیل کیفیت پایین قابل‌استفاده برای آموزش سامانه نبودند و از آنها ۹۸۳۴ قطعه شامل ۲۵۷۷۵ کلمه انگلیسی و ۲۶۰۲۹ کلمه فارسی استخراج شد.

برای ارزیابی، از سامانه ترجمه ماشینی موزس^۱ استفاده شده است. موزس یک سامانه ترجمه ماشینی آماری است که در آن می‌توانیم مدل‌های ترجمه را برای هر جفت زبان دلخواه بسازیم. کافی است مجموعه‌ای از متون آموزشی دوزبانه داشته باشیم تا پس از ساخت مدل، موزس با استفاده از یک الگوریتم جستجوی بهینه محتمل‌ترین ترجمه را ارائه کند.

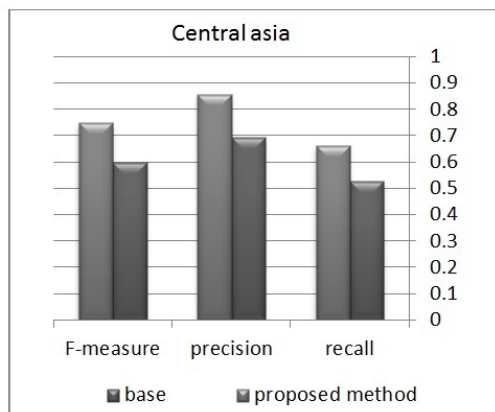
برای ارزیابی سامانه ترجمه ماشینی را یک‌بار با پیکره PEN آموزش داده و بار دیگر قطعات موازی استخراج‌شده از پیکره قابل‌مقایسه را به آن افزودیم. برای آزمون هم یک‌بار از ۱۰۰۰ جمله از پیکره PEN (مجموعه آزمون یک که البته در آموزش استفاده نشده بودند) و یک‌بار از همین تعداد جمله منتخب از پیکره khameneii.ir، Wikipedia و central asia (مجموعه آزمون ۲) استفاده کردیم. نتایج حاصل از ارزیابی بلوی سامانه عبارتند از:

¹ moses



(شکل ۴-۴): مقایسه سامانه پایه و پیشنهادی روی پیکره Khameneii.ir

در این شکل نیز مشاهده می‌شود که نتایج حاصل از روش پیشنهادی از روش پایه بهتر است. بازخوانی به‌نسبه پایین الگوریتم‌ها به‌دلیل وجود کلمات سنگین و سخنرانی‌گونه و عدم تطابق کامل حوزه متون با جدول LLR است؛ اما دقت برنامه مناسب است. تفاوت دقت الگوریتم پیشنهادی به‌دلیل استفاده از ویژگی‌های دقیق‌تر به نسبت روش پایه است.

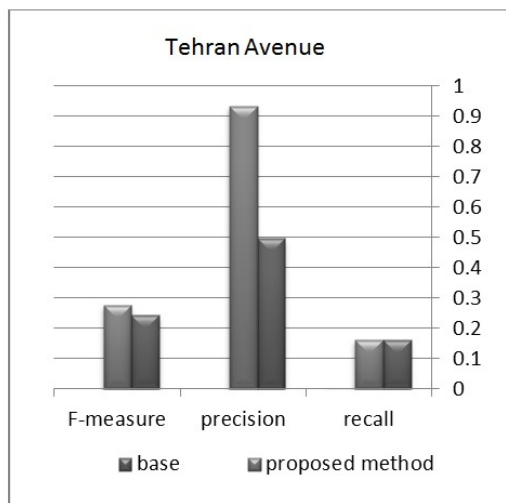


(شکل ۴-۵): مقایسه سامانه پایه و پیشنهادی روی پیکره Central Asia

همان‌طور که در نمودار مشاهده می‌شود، مشابه پیکره‌های آزمون شکل (۴-۵) دقت برنامه بالاست که دلیل آن پیش‌تر توضیح داده شد؛ اما نکته‌ای که در اینجا مشاهده می‌شود، این است که بازخوانی برنامه و در نتیجه کارایی آن نیز بسیار بالاست که این امر به‌دلیل تطابق حوزه متون

(جدول ۴-۷): مقایسه نمونه‌های ترجمه در دو سامانه ترجمه

according to law enforcement sources , five tajik soldiers were wounded and six militants were killed in the ensuing fire fight .	جمله مبدا	۱
به گفته منابع اعمال قانون , پنج سرباز تاجیکستان تن دیگر نیز زخمی شده‌اند و شش شبه‌نظامیان کشته شده در این مبارزه می‌گرفتند.	ترجمه سامانه پایه	
به گفته منابع اعمال قانون , پنج سرباز تاجیکستان زخمی شدند و شش نفر از پیکارجویان کشته‌شده در این مبارزه می‌کرد .	ترجمه سامانه با دادگان افزوده	
it was so dangerous for me and for the patients .	جمله مبدا	۲
این خطرناک بود برای من و برای بیماران سودمند است .	ترجمه سامانه پایه	
این بود بسیار خطرناک برای من و برای بیماران .	ترجمه سامانه با دادگان افزوده	



(شکل ۴-۳): مقایسه سامانه پایه و پیشنهادی روی پیکره Tehran avenue

همچنین عدم انجام ترازبندی حریصانه و در نظر گرفتن ترازبندی چندگانه که با ذات زبان فارسی مطابقت بیش‌تری دارد نیز در نتیجه بهتر مؤثر است. بازخوانی پایین برنامه همان‌طور که پیش از این هم ذکر شد، به‌دلیل عدم تطابق حوزه متون پیکره با جدول LLR است که استخراج قطعات زیاد را ناممکن ساخته است.

آن استفاده از فرهنگ لغت خاص آن موضوع پیشنهاد می‌شود.

۶- مراجع

B. Zhao and S. Vogel, "Adaptive Parallel Sentences Mining from Web Bilingual News Collection," in Proceedings of the 2002 IEEE International Conference on Data Mining. IEEE Computer Society, pp. 745-748.

C. Tillmann, "A beam-search extraction algorithm for comparable data," Proceedings of the ACL-IJCNLP 2009 Conference Short Papers, pp. 225-228, Suntec, Singapore, 4 August.

C. Tillmann, Jian-ming Xu, "A simple sentence-level extraction algorithm for comparable data," Proceedings of NAACL HLT 2009: Short Papers, pp 93-96, Boulder, Colorado.

Cerny, V. , "A thermodynamical approach to the traveling salesman problem: An efficient simulation algorithm", Journal of Optimization Theory and Applications, 1985, Vol. 45, pp. 41-51.

D. Munteanu and D. Marcu, "Extracting parallel sub-sentential fragments from non-parallel corpora," in Proceedings of the 21st International Conference on Computational Linguistics and the 44th annual meeting of the Association for Computational Linguistics. Association for Computational Linguistics, 2006 , pp. 81-88.

D. Munteanu and D. Marcu, "Improving machine translation performance by exploiting non-parallel corpora," Computational Linguistics, 2005, vol. 31, no. 4, pp. 477-504.

D. Melamed, "Bitext mapping via Pattern Recognition", Computational Linguistics , Published by MIT Press, 1999, pp 107 - 130 Maps and Alignment.

Eric Gaussier, Jean-Michel Renders, Irina Mat-veeva, Cyril Goutte, and Herve Dejean, "A geometric view on bilingual lexicon extraction from comparable corpora", In ACL 2004, pages 527-534.

Harold W. Kuhn, "The Hungarian Method for the assignment problem", Naval Research Logistics Quarterly, 1995, 2:83-97, Kuhn's original publication.

Kirkpatrick, S. , Gelatt, C. D. , and Vecchi, M. P., (1983), "Optimization by simulated annealing, Science", Vol. 220, pp. 671-680.

L.Lee, A.Aw, M.Zhang, and H.Li, "EM-based hybrid model for bilingual terminology extraction from comparable corpora", Coling 2010: Poster Volume, pp. 639-646, Beijing.

آزمون با پیکره سازنده جدول LLR است و کلمات، هم‌پوشانی بیش‌تری دارند. نکته دیگر اینکه از بخش دیگری از این پیکره در بخش تعیین وزن ویژگی‌ها استفاده شده بود. در نتیجه ساختار متنی این پیکره مناسب با وزن‌های ویژگی‌ها بوده و نتایج بسیار مطلوب است.

۵- نتیجه‌گیری

ما در این نوشتار به معرفی روشی جدید برای استخراج قطعه از پیکره قابل‌مقایسه پرداختیم. برای تحقق این امر الگوریتمی طراحی و پیاده‌سازی شد که با استفاده از جدول کلمات LLR و الگوهای رایج برای ساختار موجود قطعه‌های فارسی و انگلیسی و همچنین ویژگی‌هایی از قبیل طول قطعه، مربعی بودن بلوک قطعه و درصد کلمات هم‌ترجمه در قطعه که به‌صورت بهینه وزن‌دار شده‌اند، قطعات موازی را از متون قابل‌مقایسه استخراج می‌کند. فرآیند استخراج قطعه به این صورت است که برای هر جفت‌جمله نامزد، ماتریس امتیازات LLR متناظر جفت‌کلمات ایجاد می‌شود؛ سپس به‌صورت نزولی از نقطه دارای بیش‌ترین امتیاز، به‌عنوان نقطه شروع استخراج قطعه آغاز شده و پررنگ‌ترین نقاط یعنی نقاط دارای بالاترین امتیاز تا وقتی که از حد آستانه‌ای افت نکرده‌اند، برای استخراج قطعه، یکی پس از دیگری پردازش می‌شوند. نتایج ارزیابی‌ها نشان می‌دهد که داده‌های استخراج‌شده از یک‌سری متون قابل‌مقایسه موجود، عملکرد سامانه ترجمه ماشینی پایه را با وجود ناچیز بودن حجم دادگان افزوده‌شده به نسبت پیکره پایه آموزشی از ۰/۳۳ تا ۱/۴ واحد بلو بسته به مجموعه آزمون افزایش داده است.

همان‌طور که نتایج نشان داد، دقت برنامه بسیار بالاست؛ اما ایده‌آل بودن بازخوانی برنامه تا حد زیادی منوط به تطابق حوزه متن آزمون با پیکره سازنده جدول LLR است.

در راستای تکمیل و بهبود روش پیشنهادی، کارهای آینده‌ای از قبیل اعمال روش‌های پیش‌پردازشی دقیق‌تر مانند ریشه‌یابی و دقیق‌تر ساختن توکن‌بندی برنامه برای تشخیص تاریخ‌ها، اعداد و ... به دادگان ورودی، بهبود فرهنگ لغت ساده با استفاده از پیکره موازی بزرگ‌تر برای افزایش دقت در بخش انتخاب جملات نامزد، تکمیل جدول کلمات LLR با استفاده از داده‌های حوزه‌های مختلف و پیکره آموزشی بزرگ‌تر و امکان انتخاب موضوع متن و به تبع

Yonggang Deng, Shankar Kumar, and William Byrne, "Segmentation and alignment of parallel text for statistical machine translation". Journal of Natural Language Engineering. to appear. 2006.

Young, Steve, Corpus-Based Methods in Language and Speech Processing, Kluwer Academic Publisher, 1997.



زینب رحیمی در سال ۱۳۶۶ در شاهرود متولد شد. تحصیلات تا مقطع دیپلم را در شهر شاهرود سپری و دیپلم متوسطه خود را در سال ۱۳۸۳ دریافت کرد. وی تحصیلات خود را در مقطع کارشناسی در رشته مهندسی کامپیوتر (نرم‌افزار) در دانشگاه صنعتی امیرکبیر (۱۳۸۹) و کارشناسی ارشد را در رشته مهندسی فناوری اطلاعات (سامانه‌های چندرسانه‌ای) در دانشگاه صنعتی امیرکبیر (۱۳۹۱) به پایان رساند. ایشان هم‌اکنون دانشجوی مقطع دکترا در رشته مهندسی کامپیوتر (هوش مصنوعی) در دانشگاه شهید بهشتی تهران هستند. از موضوعات مورد علاقه ایشان می‌توان به پردازش زبان طبیعی، ترجمه ماشینی و هستان‌شناسی اشاره کرد. نشانی رایانامه ایشان عبارت است از:

rahimi.zeinab@gmail.com



محمدحسین ثمنی در سال ۱۳۶۵ در تهران متولد شد. ایشان تحصیلات خود را در مقطع کارشناسی در رشته مهندسی فناوری اطلاعات در دانشگاه صنعتی امیرکبیر (۱۳۸۸) و کارشناسی ارشد را در رشته مهندسی کامپیوتر (معماری شبکه‌های کامپیوتری) در همان دانشگاه (۱۳۹۱) به پایان رساند. ایشان هم‌اکنون پژوهش‌گر گروه امنیت زیرساخت (پژوهشکده افتا) در پژوهشگاه توسعه فناوری‌های پیشرفته خواجه نصیرالدین طوسی هستند. از موضوعات مورد علاقه ایشان می‌توان به امنیت زیرساخت، هوش مصنوعی، داده‌کاوی و پردازش داده‌های کلان اشاره کرد. نشانی رایانامه ایشان عبارت است از:

mhsamani@gmail.com

Mona Diab and Steve Finch, "A statistical word-level translation model for comparable corpora". In RIAO 2000.

Papineni, K., Roukos, S., Ward, T., & Zhu, W.-J. BLEU: a method for automatic evaluation of machine translation, RC22176 (Technical Report). IBM T.J. Watson Research Center, 2001.

Pascale Fung and Lo Yuen Yee, "An IR approach for translating new words from nonparallel, comparable texts". In ACL 1998, pages 414-420.

Pascale Fung and Percy Cheun, "Mining very non-parallel corpora: Parallel sentence and lexicon extraction via bootstrapping and EM". In EMNLP 2004, pages 57-63.

Pascale Fung and Percy Cheung. Mining Very Non-Parallel Corpora: Parallel Sentence and Lexicon Extraction via Bootstrapping and EM. In D. Lin and D. Wu, editors, Proceedings of Empirical Methods on Natural Language Processing (EMNLP'04), 2002, pp. 57-63, Barcelona, Spain.

Philipp Koehn and Kevin Knight, "Estimating word translation probabilities from unrelated mono-lingual corpora using the EM algorithm", In National Conference on Artificial Intelligence, 2004, pages 711-715

R.Barzilay, N.Elhadad "Sentence alignment for mono-lingual comparable corpora," Proceedings of the conference on Empirical methods in natural language processing, 2003, pp. 25 - 32.

R.C.Moore, "Fast and Accurate Sentence Alignment of Bilingual Corpora", Proceedings of the 5th Conference of the Association for Machine Translation in the Americas on Machine Translation, 2002, pp 135 - 144.

S. Abdul-Rauf and H. Schwenk, "Exploiting comparable corpora with ter and terp," in BUCC 09: Proceedings of the 2nd Workshop on Building and Using Comparable Corpora, 2009, pp. 46-54.

Sanjika Hewavitharana, Nguyen Bach, Qin Gao, Vamshi Ambati and Stephan Vogel, "CMU Haitian Creole-English Translation System", 2011, In Proc. WMT. Pp386-392.

W.A.Gale, K.w.church, "A Program for Aligning Sentences in Bilingual Corpora", Proceedings of the 29th annual meeting on Association for Computational Linguistics, 1991, pp 177 - 184.

X.Ma, "Champollion: a robust parallel Text sentence aligner," Proceedin of the fifth international conference on language resource and evaluation (LREC 2006), pp. 167-172.



شهرام خدیوی مدارک کارشناسی و

کارشناسی ارشد خود را در رشته

مهندسی کامپیوتر از دانشگاه صنعتی

امیرکبیر به ترتیب در سال‌های ۱۳۷۵ و

۱۳۷۸ دریافت کرده و همچنین مدرک

دکترای خود را در سال ۱۳۸۷ از دانشگاه آخن RWTH در

رشته علوم کامپیوتر دریافت کرده‌اند. ایشان در حال حاضر

عضو هیأت علمی و استادیار دانشکده مهندسی کامپیوتر

دانشگاه صنعتی امیرکبیر هستند. زمینه‌های پژوهشی مورد

علاقه ایشان پردازش زبان طبیعی، ترجمه ماشینی آماری و

یادگیری ماشین است.

نشانی رایانامه ایشان عبارت است از:

khadivi@aut.ac.ir

فصلنامه



سال ۱۳۹۴ شماره ۲ پیاپی ۲۴