

## طراحی و پیاده سازی یک سامانه ترجمه فارسی به انگلیسی

فائزه میرزائی<sup>۱</sup>، محسن بیگلری<sup>۲</sup>، احمد یوسفان<sup>۳</sup> و عماد بیات<sup>۴</sup>

دانشگاه کاشان؛ دانشکده مهندسی، گروه مهندسی کامپیوتر

[fmirzaei@grad.kashanu.ac.ir](mailto:fmirzaei@grad.kashanu.ac.ir)<sup>۱</sup>

[biglari@grad.kashanu.ac.ir](mailto:biglari@grad.kashanu.ac.ir)<sup>۲</sup>

[yoosofan@kashanu.ac.ir](mailto:yoosofan@kashanu.ac.ir)<sup>۳</sup>

[emadbayat85@yahoo.com](mailto:emadbayat85@yahoo.com)<sup>۴</sup>

### چکیده

ترجمه‌ی خودکار ماشینی از شاخه‌های فعال پردازش زبان طبیعی در سال‌های اخیر بوده است. هر سیستم پردازش زبان طبیعی نیاز به یک بخش پیش پردازش برای تبدیل ورودی به شکل مناسب دارد. این پیش پردازش می‌تواند شامل جداسازی عبارت ورودی به جملات، جملات به کلمات و اصلاح املا کلمات در صورت نیاز باشد.

در این مقاله یک سیستم ترجمه‌ی فارسی به انگلیسی با استفاده از روش مبتنی بر قوانین ارائه شده است که بخش پیش پردازش را نیز شامل می‌شود. در این سیستم بخش مهمی از پردازش‌ها به زمان طراحی سیستم منتقل شده است؛ و به موجب این انتقال، سرعت و دقت سیستم در تجزیه‌ی جملات افزایش چشمگیری پیدا کرده است. برای واژگان فارسی لغت‌نامه دهخدا به کار گرفته شده و برای ترجمه‌ی انگلیسی یک فرهنگ واژگان دوزبانه با نزدیک به ۴۵۰۰۰ لغت به کار گرفته شده است. در زمان اجرای سیستم جاری، تنها نیاز به تشخیص مواردی چون نشانه‌های جمع، زمان افعال و ضمائر متصل مفعولی می‌باشد. سیستم ارائه شده به صورت کامل پیاده‌سازی و ارزیابی شده است و نتایج بدست آمده، نشان داده است که چنین روشی کاملاً کارآمد است.

### کلمات کلیدی

ترجمه‌ی ماشینی، ترجمه‌ی فارسی به انگلیسی، لغت‌نامه دهخدا، پردازش زبان طبیعی

روش‌های مبتنی بر قوانین: این روش‌ها از مجموعه‌ای از قوانین برای ترجمه بهره می‌برند که براساس دانش زبان‌شناسی انسان نوشته شده است.

روش‌های آماری: این روش‌ها، براساس مثال‌هایی از ترجمه‌ی انجام شده توسط انسان سعی بر یادگیری ترجمه دارند.

### ۲- مشکلات ترجمه فارسی به انگلیسی

گام اول در ترجمه، تجزیه‌ی جملات زبان مبدا است؛ از دشواری‌های تجزیه‌ی جملات زبان فارسی، می‌توان به موارد زیر اشاره کرد:

- بعضی از اسم‌ها دارای چند معنی متفاوت هستند، مانند شیر.
- بطور معمول، حروف صدادار کوچک در نوشتار فارسی ظاهر نمی‌شوند که موجب ابهام در کلمات متشابه می‌شود. برای مثال، کلمه‌ی «گل» را می‌توان به دو صورت «گِل» و «گُل» تفسیر کرد.

### ۱- مقدمه

ترجمه‌ی یک متن از یک زبان طبیعی به زبانی دیگر، ترجمه‌ی ماشینی نامیده می‌شود. در ترجمه‌ی خودکار ماشینی، عملیات ترجمه توسط کامپیوتر صورت می‌گیرد. این شاخه از پردازش زبان طبیعی در ۵۰ سال اخیر، موضوع پژوهش بسیاری از محققان بوده و توجه زیادی را به خود جلب کرده است. روش‌های ارائه شده را می‌توان به صورت کلی به چند بخش تقسیم کرد:

روش‌های مبتنی بر فرهنگ واژگان: در این روش‌ها، منبع اصلی ترجمه یک فرهنگ واژگان دوزبانه است که بسیار تعیین کننده بوده و قدرت ترجمه نیز تا حدود زیادی به آن بستگی دارد.

روش‌های مبتنی بر دانش یا مثال: در این روش‌ها، با استفاده از دانش و نوشتجات دوزبانه (مثال‌ها)، تک‌تک موارد از زبان مبدا به زبان مقصد نگاشت می‌شوند.

- انعطاف و تنوع جملات فارسی: قید در هر جای جملات (غیر از آخر) می تواند قرار گیرد. بسته به شرایط، ترتیب مفعول و متمم جابه جا می شود. اگر مفعول با «را» بیاید، قبل از متمم قرار می گیرد و اگر با «ی» بیاید، هم قبل و هم بعد از متمم می تواند قرار گیرد؛ و در صورتی که هیچکدام از این دو را نداشته باشد، معمولاً بعد از متمم قرار می گیرد. در زبان عامیانه گاهی فعل در ابتدای جمله ظاهر می شود.
- حرف اضافه برای ارتباط دادن دو کلمه (اسم و صفت)، معمولاً در نوشتار فارسی آورده نمی شود. مانند «گل قرمز»: «گل قرمز».
- وجود استثنائات زیاد در قوانین فارسی: برای مثال، علامت جمع فارسی «ان» اگر به کلمه ای که با «ه» پایان می یابد، اضافه شود، به «گان» تبدیل می شود.

### ۳- مروری با کارهای انجام شده

به صورت کلی، پروژه های انگشت شماری در رابطه با ترجمه ی فارسی به انگلیسی و برعکس صورت گرفته است؛ که اغلب آن ها نیز از روش های مبتنی بر قوانین استفاده کرده اند. پروژه ی شیراز [1] یکی از سیستم های نخستین در این رابطه است که با استفاده از یک ساختار از پیش تعریف شده و یک فرهنگ واژگان تقریباً ۵۰۰۰۰ کلمه ای پیاده سازی شده است. این سیستم از یک تجزیه کننده ی نحوی و تحلیل گر ریخت شناسی<sup>۱</sup> بهره برده است. مراجع [2-4] نیز با استفاده از گرامر درخت مجاورت<sup>۲</sup> یا TAG یک سیستم ترجمه ی انگلیسی به فارسی با استفاده از روش مبتنی بر قوانین ارائه کردند؛ سپس نسخه ی توسعه یافته ای از این سیستم را با استفاده از درخت های تصمیم آموزش دیده معرفی کردند. مترجم های تجاری همچون پارس و آریا نیز موجودند که دارای دقت بالایی نیستند. در [5] یک مترجم دوطرفه ی فارسی به انگلیسی و انگلیسی به فارسی به نام PEnTrans ارائه شده است؛ این مترجم از دو بخش اصلی با نام های PEnT1 برای ترجمه ی انگلیسی به فارسی و PEnT2 برای ترجمه ی فارسی به انگلیسی تشکیل شده است؛ در این مرجع، پس از معرفی سیستم مربوطه، مقایسه ای بین کارکرد آن و سیستم های موجود نیز صورت گرفته است تا برتری خود را نشان دهد.

### ۴- روش پیشنهادی

سیستم ارائه شده در این مقاله، از ترکیبی از روش مبتنی بر قوانین و فرهنگ واژگان بهره می برد. این سیستم از دو بخش اصلی تشکیل شده است:

۱. تجزیه ی ورودی فارسی

۲. تبدیل به معادل انگلیسی

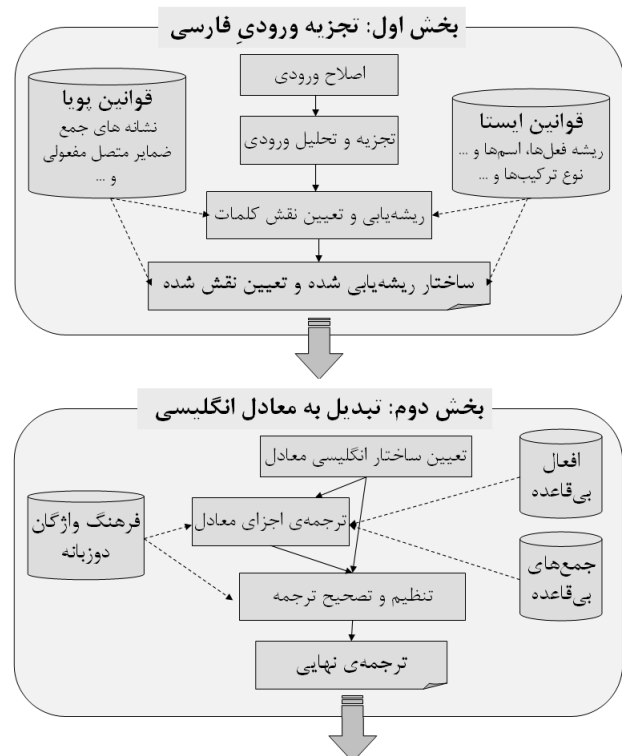
بخش اول بسیار حائز اهمیت است زیرا از جهات مختلفی تعیین کننده است؛ هر اشتباهی در این بخش، منجر به ترجمه ی نادرست در بخش بعد خواهد شد. از طرفی، سرعت کل سیستم به صورت عمده توسط این بخش تعیین می شود. مشکل اغلب روش های ارائه شده در تحلیل و تجزیه ی جملات فارسی، سرعت پردازش است؛ این روش ها در زمان اجرا، به تحلیل کلمات در سطح های مختلف نحوی و معنایی می پردازند که موجب کاهش عمده ی سرعت می شود؛ در صورتی که بخش بزرگی از این پردازش ها تکراری است و قابل انجام در زمان طراحی: مانند ریشه یابی کلمات. از این رو، سیستم ارائه شده، سعی بر افزایش چشم گیر سرعت با انتقال حداکثر مقدار ممکن از پردازش ها به زمان طراحی دارد. این بهینه سازی به بخش اول سیستم اعمال شده است و در بخش مربوطه به تفصیل توضیح داده شده است.

سیستم در گام اول یک فرهنگ لغت غنی تشکیل می دهد. سپس قوانین ساخت کلمات را بر روی این فرهنگ لغت اعمال کرده و به تجزیه ی کلمات در سطح های مختلفی می پردازد؛ در نهایت، حاصل پردازش ها با یک ساختار مشخص در پایگاه داده قرار داده می شود. با استفاده از این ساختار، هر کلمه می تواند بدون هیچ پردازش اضافی، به اجزای تشکیل دهنده ی خود تجزیه شود. در زمان اجرا، سیستم با استفاده از پایگاه داده ی ایجاد شده، اقدام به تجزیه ی جملات ورودی کرده و نقش بخش های مختلف را تعیین می کند. حاصل این تجزیه، یک ساختار ریشه یابی شده و تعیین نقش شده است که برای بخش دوم سیستم ارسال می شود. بخش دوم سیستم که مسئول ترجمه ی انگلیسی است، ابتدا بر اساس ساختار جمله ی فارسی، ساختار انگلیسی معادل را محاسبه می کند؛ سپس با استفاده از فرهنگ واژگان دوزبان، کلمات مورد نظر را ترجمه کرده و بر اساس ساختار مورد نظر، در قالب مشخصی قرار می دهد؛ در این مسیر، افعال بی قاعده و اسم هایی که دارای جمع های بی قاعده هستند نیز در نظر گرفته می شوند. در نهایت پس از تنظیم و تصحیح احتمالی، خروجی نهایی تولید می شود. شکل ۱، سیستم ارائه شده را نمایش می دهد.

### ۴-۱- بخش اول: تجزیه ی جملات فارسی

تشخیص درست و شناسایی کامل کلمات و عبارات، پیش نیاز یک ترجمه ی صحیح است. گام اول در شناسایی کلمات، ریشه یابی آن ها است؛ کلمات در زبان فارسی به صورت

- کلی به انواع فعل، اسم، صفت، قید، ضمیر و حرف اضافه تقسیم می‌شوند [۶].
- در سیستم ارائه شده از یک پایگاه داده‌ای توسعه یافته استفاده می‌شود؛ در این پایگاه داده برای دسته‌بندی کلمات و قرار دادن آن‌ها در پایگاه داده، کلمات همراه با نوع، بخش‌های تشکیل دهنده و اطلاعات اضافه‌تری قرار گرفته‌اند.



شکل (۱): ساختار سیستم ارائه شده

#### ۴-۱-۱- ایجاد پایگاه داده‌ی ایستا

برای افزایش سرعت و کاهش ضریب خطا در ریشه‌یابی و تجزیه‌ی کلمات، قوانین ساخت کلمات به صورت ایستا پردازش و در پایگاه داده‌ای قرار داده می‌شود.

مرجع اصلی استفاده شده برای کلمات فارسی استفاده شده، فرهنگ دهخدا بوده است که در [7] درباره‌ی چگونگی ساخت آن توضیح داده شده است. کلمات براساس نوعشان در گروه‌های زیر از فرهنگ دهخدا استخراج شده‌اند:

- مصدر فعل
- بن ماضی
- بن مضارع
- اسم
- صفت
- قید
- ضمیر

کلمات استخراج شده، بازبینی و کلمات کم‌کاربرد و خاص هرس شده‌اند. این عمل موجب افزایش سرعت شده و از طرفی در مراحل بعدی کار، از تولید نتایج ناخواسته و دور از انتظار جلوگیری می‌کند. در نهایت و پس از پیش‌پردازش‌های موردنیاز، کلمات براساس گروه مربوطه در پایگاه داده قرار داده می‌شوند. در ساخت پایگاه داده، دو سطح قراردادی در نظر گرفته می‌شود؛ در سطح قراردادی اول، قوانینی که دارای درصد اطمینان بالاتر هستند، بر روی کلمات موجود در فرهنگ لغت اعمال می‌شوند و در سطح قراردادی دوم، قوانینی که دارای درصد اطمینان پایین‌تر هستند، اعمال می‌شوند (جدول ۱). این روش موجب حذف بسیاری از ترکیبات نادرست می‌شود. جدول ۲ و ۳ چند نمونه از قوانین سطح قراردادی اول و دوم را نمایش می‌دهند.

#### جدول (۱): قوانین مربوط به سطح‌های قراردادی اول و دوم

بن افعال + وند (پس/پیش)	سطح قراردادی اول
اسم، صفت، ضمیر، قید + وند (پس/پیش)	سطح قراردادی دوم

#### جدول (۲): چند نمونه از قوانین سطح قراردادی اول برای اسم‌ها

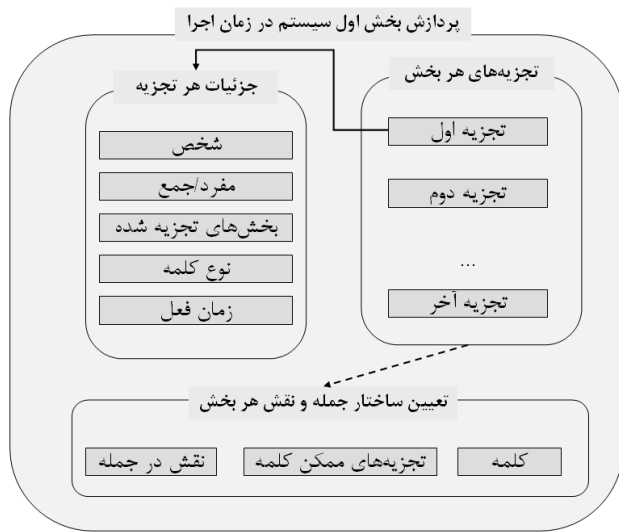
نتیجه	بخش سوم	بخش دوم	بخش اول
اسم مصدر	-	ش	بن مضارع
اسم فاعلی	-	نده	بن مضارع
اسم مرکب	-	گاه	بن مضارع
اسم مصدر/مرکب	بن مضارع	و	بن مضارع

#### جدول (۳): چند نمونه از قوانین سطح قراردادی دوم برای اسم‌ها

نتیجه	بخش سوم	بخش دوم	بخش اول
اسم	-	چه	اسم
اسم مرکب	-	اسم	اسم
اسم	-	بن مضارع	اسم
اسم مکان	-	بن مضارع	صفت

برای هر کلمه در پایگاه داده ساختار زیر در نظر گرفته شده است:

- یک کلید یکتا
- نوع کلمه
- سطح قراردادی
- اشاره‌گر به بخش اول کلمه (در صورت وجود)
- اشاره‌گر به بخش دوم کلمه (در صورت وجود)
- اشاره‌گر به بخش سوم کلمه (در صورت وجود)
- اطلاعات مرتبط با کلمه: این بخش برای هر نوع کلمه، می‌تواند کاربرد منحصری داشته باشد. برای مثال، برای افعال، می‌تواند مصدر/بن ماضی/بن مضارع آن را مشخص



شکل (۳): روند عملیات سیستم برای پردازش یک جمله‌ی ورودی

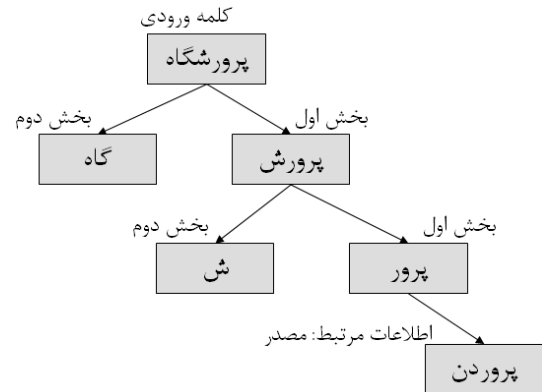
در سیستم پیاده‌سازی، جملات خبری ساده مورد بررسی قرار گرفتند. هر جمله‌ی خبری به صورت کلی به دو بخش گروه نهادی و گروه گزاره‌ای قابل تقسیم است. گروه نهادی از اولین کلمه‌ی جمله آغاز و با حرف‌های ربطی همچون «و» و «یا» گسترش می‌یابد. باقی جمله نیز گروه گزاره‌ای خواهد بود. در صورتی که جمله دارای فعل اسنادی باشد، نقش بخش اول «نهاد» و در غیراینصورت «فاعل» خواهد بود. برای تعیین نقش هر کلمه در بخش دوم (گروه گزاره‌ای)، از اطلاعات موجود در جمله و اطراف هر کلمه استفاده می‌شود؛ برای مثال، مفعول بعد از نشانه‌ی مفعولی ظاهر می‌شود. نتیجه‌ی اجرای جمله‌ی ورودی «بچه و پدرش به مدرسه می‌روند» در شکل ۴ نمایش داده شده است. ساختارهای پشتیبانی شده در جدول ۴ فهرست شده‌اند.

جدول (۴): ساختارهای قابل شناسایی در تجزیه جملات

مثال	ساختار
باران می‌بارد	فاعل + فعل
کارگر آجرها را حمل می‌کند	فاعل + مفعول + فعل
مادر کودکش را به مدرسه می‌برد	فاعل + مفعول + متمم + فعل
هوا سرد است	نهاد + مسند + فعل اسنادی
دانشجو به دانشگاه می‌رود	فاعل + متمم + فعل
راننده با بوق مسیرش را باز کرد	فاعل + متمم + مفعول + فعل
هوا مانند زمستان سرد است	نهاد + متمم + مسند + فعل اسنادی

کند. برای اسم‌ها، می‌تواند اسامی هم‌معنی آن را مشخص کند. برای صفت‌ها، می‌تواند شامل صفات هم‌معنی/متضاد آن باشد.

پایگاه داده‌ی جاری به شکلی طراحی شده است که با مراجعه به هر کلمه، می‌توان به بالاترین سطح از تجزیه‌ی آن دسترسی پیدا کرد. هر کلمه دارای یک کلید یکتا است و از یک یا چند بخش تشکیل شده است؛ هر بخش خود به کلمه‌ای دیگر اشاره می‌کند که می‌توان برای تجزیه‌ی بیشتر به آن مراجعه کرد. برای مثال کلمه‌ی «پرورشگاه» را در نظر بگیرید؛ که از دو بخش «پرورش» و «گاه» تشکیل شده است؛ حال اگر به جزئیات بیشتری نیاز باشد، می‌توان به کلمه‌ی «پرورش» مراجعه کرد که خود از دو بخش «پرور» و «ش» تشکیل شده است؛ با مراجعه به کلمه‌ی «پرور» می‌توان پی برد که این کلمه بن مضارع مصدر «پروردن» است.



شکل (۲): روند تجزیه‌ی یک کلمه تا آخرین سطح با استفاده از پایگاه داده‌ی طراحی شده

#### ۴-۱-۲- پردازش پویا

منظور از پردازش پویا، اعمال سیستم در زمان اجراست؛ که شامل تجزیه‌ی جملات به بخش‌های کوچکتر و یافتن نقش هر بخش می‌باشد. هر کلمه در جمله ممکن است به چند طریق قابل تجزیه باشد، بنابراین بعد از جداسازی کلمه‌های موجود در جمله، تجزیه‌های ممکن هر بخش استخراج می‌شود؛ سپس برای هر تجزیه، اطلاعات مرتبط با آن محاسبه می‌شود؛ در نهایت مجموع این اطلاعات برای زیربخش بعدی که وظیفه‌ی تعیین نقش هر بخش از جمله را دارد، ارسال می‌شود. شکل ۳ روند عملیات را نمایش می‌دهد.

گروه نهادی (فاعل)

کلمه	سایر اطلاعات
بچه	اسم، مفرد
و	حرف اضافه، مفرد
پدرش	اسم، [پدر، ش]، سوم شخص مفرد

گروه گزاره‌ای

کلمه	نقش	سایر اطلاعات
به	نشانه متمم	حرف اضافه، مفرد
مدرسه	متمم	اسم، مفرد
می‌روند	فعل	فعل، [می، رو، ند]، سوم شخص جمع، مضارع اخباری، از مصدر رفتن

شکل (۴): نتیجه‌ی تجزیه یک نمونه جمله‌ی ورودی

## ۵- بخش دوم: ترجمه‌ی انگلیسی

فرهنگ واژگان دوزبانه‌ی کامل یکی از ضرورت‌های مهم در ترجمه‌ی انگلیسی است. در سیستم جاری از یک فرهنگ واژگان ۴۵۰۰۰ کلمه‌ای استفاده شده است. روش کار بخش دوم سیستم به این شکل است که برای هرکدام از هفت ساختار استخراج شده توسط بخش اول، ساختار معادل انگلیسی را در نظر می‌گیرد. در صورتی که ساختار ورودی یکی از این هفت حالت نباشد، از ترجمه‌ی بخش به بخش استفاده شده است؛ در این حالت بهینه‌سازی‌هایی بر روی ورودی و خروجی نیز اعمال می‌شود؛ هرچند در نهایت خروجی به دقت حالت‌های قبیل نیست. اسامی خاص با استفاده از حروف معادل انگلیسی جایگزین می‌شوند؛ هرچند به دلیل عدم وجود حرکت در اسامی خاص فارسی، ممکن است معادل‌سازی به صورت دقیق انجام نشود.

مراحل کلی بخش دوم از سیستم به صورت زیر است:

۱. تعیین ساختار انگلیسی معادل با ساختار فارسی ورودی

۲. ترجمه‌ی اجزای جمله‌ی ورودی به انگلیسی

۳. ترکیب اجزای ترجمه شده براساس ساختار مربوطه

جدول ۵ مثالی از چند جمله‌ی ورودی و معادل ترجمه‌ی شده‌ی آن‌ها را توسط سیستم ارائه شده نشان می‌دهد.

جدول (۵): چند نمونه از ترجمه‌های انجام شده توسط سیستم جاری

جمله‌ی ورودی	ترجمه
آریا آمد	Aria came
او با من آمد	He/She came with me
ناهید پولش را به من بخشید	Nahid spared her/his deposit to me

## ۶- نتیجه‌گیری

برای بررسی کارایی و سرعت سیستم ارائه شده، آن را توسط زبان برنامه‌نویسی پایتون<sup>۳</sup> و پایگاه داده‌ای SQLite که

هر دو متن‌باز و رایگان هستند، پیاده‌سازی و ارزیابی کردیم [8]. نتایج در محدوده‌ی ساختارهای در نظر گرفته شده بسیار مناسب و دقیق بود. برای محاسبه‌ی یک مقیاس تقریبی از دقت سیستم، ۱۰۰ جمله در قالب هفت ساختار اشاره شده به عنوان ورودی به سیستم داده شد و بعد از بررسی دستی خروجی‌ها، نتایج زیر حاصل شد:

- ۸۳ جمله به صورت درست ترجمه شده بود.
- ۱۰ جمله به صورت تقریبی درست ترجمه شده بود.
- در ۷ جمله‌ی باقیمانده، از واژگان مناسب استفاده نشده بود از مزایای این سیستم نسبت به سایر سیستم‌های موجود، می‌توان موارد زیر را نام برد:
- سرعت بسیار مناسب در تجزیه‌ی جملات در زمان اجرا
- دقت بالای تجزیه‌ی کلمات موجود در فرهنگ واژگان و از مشکلات آن نیز می‌توان به موارد زیر اشاره کرد:
- در نظر نگرفتن همه‌ی ساختارهای موجود در جملات فارسی؛ مخصوصات جملات غیررسمی.
- عدم اعمال جنسیت (مرد/زن) در ترجمه‌ها.

## مراجع

- [1] Amtrup, J.W., et al., Persian-English machine translation: An overview of the Shiraz project. Memoranda in Computer and Cognitive Science MCCA-00-319, NMSU, CRL, 2000.
- [2] Feili, H. and G. Ghassem-Sani. An application of lexicalized grammars in English-Persian translation. 2004. Citeseer.
- [3] Feili, H. and G. Ghassem-Sani. Using Tree Adjoining Grammar in English-Persian Translation. 2004.
- [4] Pilevar, M.T. and H. Faili. Persiansmt: A first attempt to english-persian statistical machine translation. 2010. JADT.
- [5] Saedi, C., Y. Motazadi, and M. Shamsfard. Automatic translation between English and Persian texts. 2009.
- [6] گیوی، ا. . دستور زبان فارسی ۱. ویرایش سوم. ۱۳۸۷، تهران: موسسه فرهنگی فاطمی.
- [7] یوسفان، ا.، جعفری، خ. و بیگری، م. "تبدیل خودکار کلمه‌های لغت نامه دهخدا به قالب آوایی IPA"، دومین کنفرانس ملی مهندسی برق ایران. دومین کنفرانس ملی مهندسی برق ایران. دانشگاه آزاد اسلامی واحد نجف آباد، ایران: اسفند ۱۳۸۸
- [8] بیات، ع. . مترجم فارسی به انگلیسی جملات. ۱۳۸۹. دانشگاه کاشان: دانشکده مهندسی- گروه کامپیوتر.
- [9] میرزائی، ف. ، بیگری، م. . سیستم ترجمه جملات فارسی به انگلیسی. ۱۳۸۹. دانشگاه کاشان: دانشکده مهندسی- گروه کامپیوتر.



نخستین کنفرانس بین المللی پردازش خط و زبان فارسی

۱۵ و ۱۶ شهریور ۱۳۹۱

دانشگاه سمنان - دانشکده مهندسی برق و کامپیوتر

زیر نویس ها

---

<sup>1</sup> Morphological

<sup>2</sup> Tree Adjoining Grammar

<sup>3</sup> Python