



ارائه یک سامانه ترجمه ماشینی ترکیبی بر پایه رمزگشای یک‌نوا

حسین خاتمی^{*}، حکیمه فدائی و هشام فیلی

دانشکده مهندسی برق و کامپیوتر، پردیس دانشکده‌های فنی، دانشگاه تهران، تهران، ایران

چکیده

در این مقاله یک مترجم خودکار متون انگلیسی به فارسی با استفاده از معماری ترکیبی قاعده‌مند و آماری ارائه شده است. این معماری ترکیبی به منظور بهبود نتایج هر دو مترجم، خروجی مترجم ماشینی قاعده‌مند و آماری را ترکیب کرده و سعی می‌کند یک خروجی برتر از هر دو سامانه ایجاد کند. در این راستا از یک رمزگشای یک‌نوا با پیچیدگی زمانی چندجمله‌ای استفاده می‌شود. مترجم‌های ماشینی قاعده‌مند عمل ترجمه را بر اساس مجموعه‌ای از قواعد زبانی انجام می‌دهند. به طور معمول نتایج آنها از نظر ترتیب کلمات و ساختار نحوی، کیفیت بهتری نسبت به نتایج مترجم‌های آماری دارند؛ ولی عملکرد این مترجم‌ها در زمینه انتخاب لغات مناسب و روانی ترجمه، ضعیف‌تر از مترجم‌های ماشینی آماری است. از این‌رو در این معماری، ترجمه اولیه به وسیله مترجم ماشینی قاعده‌مند صورت می‌گیرد؛ سپس با استفاده از مترجم ماشینی آماری ترجمه آن بهبود داده می‌شود. به این منظور، ترتیب واژگان در ترجمه نهایی بر اساس ترجمه مترجم ماشینی قاعده‌مند صورت می‌گیرد؛ سپس عمل ترجمه و انتخاب لغات توسط رمزگشای یک‌نوا، با در نظر گرفتن ترجمه‌های نامزدهای ارائه‌شده توسط مترجم قاعده‌مند و آماری و همچنین با استفاده از مدل زبانی، انجام می‌شود. آزمایش‌های انجام‌شده نشان می‌دهند که کیفیت نتایج به دست آمده از معماری ترکیبی در معیار بلو، به طور تقریبی پنج واحد بهتر از نتایج مترجم ماشینی قاعده‌مند است. همچنین کیفیت این نتایج نسبت به نتایج مترجم ماشینی آماری در معیار بلو، یک واحد بهتر است.

واژگان کلیدی: مترجم ماشین، معماری ترکیبی، رمزگشای یک‌نوا، ترتیب کلمات ترجمه، انتخاب لغات

A Hybrid Machine Translation System Based on a Monotone Decoder

Hosein Khatami^{*}, Hakimeh Fadaei & Hesham Faili

School of Electrical and Computer Engineering, College of Engineering, University of Tehran, Tehran, Iran

Abstract

In this paper, a hybrid Machine Translation (MT) system is proposed by combining the result of a rule-based machine translation (RBMT) system with a statistical approach. The RBMT uses a set of linguistic rules for translation, which leads to better translation results in terms of word ordering and syntactic structure. On the other hand, SMT works better in lexical choice. Therefore, in our system, an initial translation is generated using RBMT. Then the proper lexical for the resulted sentence is chosen by using a decoder algorithm which is inspired by SMT architecture.

In the pure SMT approach, decoder is responsible for selecting proper final lexical during the translation procedure. Normally this method deals with lexical choice as well as reordering and required exponential order in time complexity. By fixing the word order in the output, a polynomial version of this method, named monotone decoding, is used in this paper. Monotone decoder algorithm selects the best lexical from a candidate list by maximizing the language model of resulted sentence. The candidate list is gathered from the outputs of both pure RBMT and pure SMT systems.

^{*} Corresponding author

^{*} نویسنده عهده‌دار مکاتبات

The experiments of proposed hybrid method on English-Persian language pair show significant improvements over both RBMT and SMT results. The results show that the proposed hybrid method gains an improvement of almost +5 units over RBMT and about one unit over SMT in BLEU score.

Keywords: Machine translation, hybrid architecture, monotone decoder, translation reordering, lexical choice

ساخت‌واژی بین عبارات^{۱۱}، کمتر در ترجمه آنها رعایت می‌شود [5]. علاوه بر این، مشکل دیگری که در مترجم‌های آماری وجود دارد، هزینه زیاد رمزگشای^{۱۲} آنها است. این نوع مترجم‌های ماشینی با وجود هرس بخشی از فضای جستجوی خود، در یک فضای به‌نسبه بزرگ، بخش زیادی از حالات ممکن برای ترجمه جمله مورد نظر را بررسی می‌کنند، که این عمل باعث افزایش هزینه ترجمه آنها می‌شود.

همان‌گونه که مشخص است، مترجم‌های ماشینی قاعده‌مند و آماری دارای نقاط قوت و ضعفی هستند که به‌نحوی تکمیل‌کننده یکدیگرند. از این‌رو در سال‌های اخیر توجه به معماری‌های ترکیبی^{۱۳} که با ارائه روش‌هایی، این دو نوع مترجم ماشینی را با یکدیگر ترکیب می‌کنند، بیشتر شده است. در این روش‌ها سعی می‌شود که با استفاده از نقاط قوت یک نوع مترجم، نقاط ضعف نوع دیگر پوشش داده شود تا معماری جدیدی با بهره‌مندی از نقاط قوت هر دو نوع مترجم ماشینی ایجاد شود. معماری‌های ترکیبی عمل ترجمه را به‌صورت موازی یا متوالی انجام می‌دهند. در روش‌های موازی، خروجی چندین نوع مترجم ماشینی با یکدیگر ترکیب می‌شوند تا ترجمه نهایی سامانه حاصل شود. در روش‌های متوالی نیز، عمل ترجمه بر پایه یک نوع مترجم ماشینی است؛ و نوع دیگر وظیفه تکمیل و بهبود ترجمه اولیه را دارد. در روش‌هایی که ترجمه اصلی بر پایه مترجم ماشینی قاعده‌مند است، به‌نحوی سعی می‌شود تا که با استفاده از مترجم آماری، نقطه ضعف انتخاب لغات مناسب و روانی ترجمه مترجم قاعده‌مند رفع شود. همچنین در روش‌هایی که ترجمه بر پایه مترجم آماری است، سعی می‌شود با استفاده از مترجم قاعده‌مند ترجمه‌ها از نظر نحوی و ترتیب واژگان بهبود داده شوند [13].

در این مقاله، به ارائه یک معماری ترکیبی متوالی می‌پردازیم که به منظور بهره‌مندی از نقاط قوت مترجم قاعده‌مند، عمل ترجمه را بر پایه این نوع مترجم انجام می‌دهد و سپس با استفاده از یک رمزگشای یکنوا^{۱۴}، ترجمه قاعده‌مند را به‌وسیله مترجم آماری بهبود می‌بخشد. رمزگشای یکنوا،

۱- مقدمه

ترجمه ماشینی^۱ یکی از مهم‌ترین و پرکاربردترین شاخه‌های پردازش متن است، در این زمینه پژوهش‌های زیادی در راستای بهبود عملکرد آنها صورت گرفته است. از پرکاربردترین و رایج‌ترین کارهای صورت‌گرفته در این حوزه، به روش‌های قاعده‌مند^۲ و آماری^۳ می‌توان اشاره کرد، که هر یک از آنها دارای نقاط قوت و ضعف هستند. یکی از نخستین رویکردها در حوزه مترجم‌های ماشینی، روش قاعده‌مند است. در این روش عمل ترجمه بر پایه مجموعه‌ای از قواعد زبان که می‌تواند توسط متخصصان زبان تولید شود، صورت می‌گیرد. این قواعد زبانی بر اساس زبان مبدا و مقصد، قواعد انتقال نحوی و لغوی از یک زبان به زبان دیگر را پیاده‌سازی می‌کنند. این نوع مترجم‌ها به‌دلیل استفاده از این قواعد و اطلاعات زبان‌شناسی، پرهزینه ولی از نظر تطابق فعل و فاعل، ترجمه مناسب شناسه فعل، خصوصیات ساخت‌واژی^۴ و به‌طور کلی از نظر ساختار نحوی^۵ عملکرد خوبی دارند. در مقابل، این نوع مترجم‌ها ترجمه روان و سلیسی ندارند و در انتخاب لغات^۶ مناسب، یک‌پارچگی واژگان و نزدیکی آنها به لغاتی که انسان برای ترجمه استفاده می‌کند، ضعیف هستند [5].

یکی دیگر از متداول‌ترین روش‌ها در حوزه ترجمه ماشینی، روش‌های آماری مبتنی بر عبارت^۷ است. این نوع مترجم‌ها بر اساس مجموعه‌ای از پیکره موازی^۸ که می‌تواند توسط ترجمه انسانی تولید شده باشد، عمل ترجمه را انجام می‌دهند. مترجم‌های آماری به‌دلیل استفاده از واژگان و عبارات استخراج‌شده از این پیکره و همچنین استفاده از مدل زبانی^۹، نتایج بهتری نسبت به مترجم‌های قاعده‌مند در زمینه انتخاب لغات مناسب و نزدیکی به ترجمه انسان دارند. این مترجم‌ها از نظر جابه‌جایی‌های^{۱۰} نزدیک و ترجمه‌های در سطح عبارت، عملکرد خوبی دارند؛ ولی وابستگی‌های نحوی و

¹ Machine translation

² Rule-based

³ statistical

⁴ morphology

⁵ syntax

⁶ Lexical choice

⁷ Phrase base

⁸ Parallel corpus

⁹ Language model

¹⁰ reordering

¹¹ phrase

¹² decoder

¹³ Hybrid architect

¹⁴ Monotone decoder

۲- کارهای مرتبط

معماری‌های ترکیبی با دو روش متوالی و موازی، کار بهبود مترجم‌های ماشینی را انجام می‌دهند. هرچند هر دو روش در راستای بهبود عملکرد مترجم‌های ماشینی موفق بوده‌اند، ولی ایده اصلی معماری ترکیبی تا زمانی که تعامل و ارتباطی بین مترجم‌های ماشینی مختلف برقرار نباشد، محقق نمی‌شود [13]. روش‌های موازی، خروجی سامانه‌های مختلف را با یکدیگر ترکیب می‌کنند؛ به گونه‌ای که جمله ورودی به صورت موازی، توسط چندین مترجم ماشینی ترجمه شده و سپس بخش‌های مختلف خروجی آنها به عنوان ترجمه نهایی، در کنار یکدیگر قرار می‌گیرند؛ از این رو به این روش، ترکیب سامانه‌ها نیز گفته می‌شود. پارک و همکاران [16] با ارائه یک طبقه‌بند^۵، خروجی سه مترجم آماری و دو مترجم قاعده‌مند را با یکدیگر ترکیب کرده‌اند. مقاله یادشده با استفاده از مجموعه‌ای از ویژگی‌ها^۶ که بر اساس اطلاعات زبانی و نحوی زبان‌های مبدا و مقصد هستند، یک طبقه‌بند ارائه داده است که توانایی انتخاب بهترین ترجمه برای یک عبارت را دارد. ترجمه‌های هر عبارت که توسط مترجم‌های ماشینی تولید شده‌اند، به عنوان ورودی به طبقه‌بند داده شده و خروجی آن، بهترین ترجمه برای آن عبارت است. با این ایده بخش‌های مختلف ترجمه‌های مترجم‌های ماشینی متفاوت، با یکدیگر ترکیب می‌شوند تا ترجمه نهایی سامانه ایجاد شود.

ایده ترکیب خروجی مترجم‌های ماشینی، به طور معمول با ساخته شدن شبکه ترجمه^۷ همراه است. در شبکه ترجمه هر گره به عنوان یک عبارت و هر یال به عنوان یک ترجمه نامزد با احتمال ترجمه‌اش است، هدف یافتن بهترین مسیر با بیشترین احتمال از گره نخست به آخر است. ماچری و همکارانش [14] در مقاله خود نیز با در نظر گرفتن خروجی یک مترجم ماشینی به عنوان ترجمه اصلی، با استفاده از هم‌ترازی^۸ بین جمله مبدا و ترجمه هر مترجم، ترتیب واژگان خروجی مترجم‌های ماشینی را بر اساس ترتیب آنها در ترجمه اصلی تغییر داده و بار دیگر با استفاده از هم‌ترازی بین جمله مبدا و ترجمه مترجم‌های ماشینی، شبکه ترجمه را ایجاد می‌کند. در این شبکه، عبارات مختلف جمله مبدا به عنوان گره و ترجمه‌های نامزد آن که از مترجم‌های ماشینی مختلف تولید شده‌اند، به عنوان یال خروجی از آن گره شناخته می‌شوند. مقاله معرفی شده با توجه به احتمالی که مترجم‌های

عمل جستجو در فضای حالات کل فرضیه‌های ترجمه ماشینی را که در آن امکان جابه‌جایی واژگان مبدا در نظر گرفته نمی‌شود، انجام می‌دهد. درواقع اگر فرض شود که ترتیب عبارات جمله مبدا با مقصد یکسان باشد، از این رمزگشا می‌توان بهره برد. این الگوریتم بر پایه برنامه‌نویسی پویا^۱ تعریف می‌شود و دارای هزینه زمانی $O(n^3)$ است که n تعداد عبارات جمله مبدا است.

در این سامانه ابتدا جمله ورودی به وسیله یک مترجم قاعده‌مند ترجمه شده و سپس از ترتیب واژگان آن برای ترجمه نهایی استفاده می‌شود. علاوه بر این، برای ترجمه عبارات جمله توسط رمزگشای یک‌نوا، از ترجمه‌های نامزد^۲ مترجم قاعده‌مند و آماری، با به کارگیری مدل زبانی، استفاده می‌شود. در نتیجه ترجمه نهایی از نظر ساختار نحوی و روانی ترجمه، دارای کیفیت بهتری در معیارهای مربوطه است. همان‌طور که گفته شد در این معماری، برای عمل ترجمه از رمزگشای یک‌نوا با هزینه ترجمه چندجمله‌ای^۳ استفاده شده است. همچنین به دلیل اینکه ابتدا جمله توسط یک مترجم قاعده‌مند با هزینه ترجمه چندجمله‌ای، ترجمه می‌شود، هزینه ترجمه نهایی معماری ترکیبی، چندجمله‌ای است که دارای هزینه کمتری نسبت به مترجم‌های آماری است.

در این مقاله آزمایش‌ها بر روی جفت زبان فارسی-انگلیسی انجام شده است که به ترتیب زبان‌های مبدا و مقصد هستند و نتایج نشان‌دهنده بهبود عملکرد نسبت به هر دو مترجم قاعده‌مند و آماری است. نتایج معماری ترکیبی در مقایسه با مترجم قاعده‌مند به طور تقریبی دارای افزایش پنج واحدی در معیار بلو^۴ است که این بهبود در مقایسه با مترجم آماری دارای افزایش به طور تقریبی یک واحدی در این معیار است.

ادامه مقاله به این ترتیب ارائه شده است: بخش ۲ به مرور کارهای مرتبط پیشین پرداخته است. در بخش ۳ ابتدا به معرفی معماری ترکیبی می‌پردازیم؛ سپس به توضیح رمزگشای یک‌نوا پرداخته و در ادامه نحوه احتمال‌دهی به ترجمه‌های مترجم ماشینی قاعده‌مند خود و نحوه ترکیب و انتخاب بهترین ترجمه‌ها از میان ترجمه‌های نامزد مترجم قاعده‌مند و آماری را شرح می‌دهیم. در بخش ۴ به معرفی سامانه‌ها، آزمایش‌ها و نتایج آنها پرداخته و در بخش آخر نیز به نتیجه‌گیری و کارهایی که می‌توان در آینده انجام داد، اشاره‌ای می‌شود.

⁵ classifier⁶ features⁷ Translation network⁸ alignment¹ Dynamic programming² Candidate translation³ Polynomial complexity⁴ BLEU measure

ماشینی مختلف به ترجمه خود داده‌اند، برای هر عبارت بهترین ترجمه را انتخاب می‌کند و در نهایت بهترین مسیر با بیشترین احتمال ترجمه در این شبکه یافت می‌شود.

تفاوت اصلی کارهای معرفی‌شده با مترجم ماشینی ارائه‌شده در این مقاله در روش کار آنها است. درواقع کارهای معرفی‌شده از روش موازی برای عمل ترجمه استفاده می‌کنند؛ درحالی‌که روش ارائه‌شده در این مقاله بر پایه روش‌های متوالی است. در روش متوالی از الگوریتمی برای انتخاب بهترین ترجمه از میان خروجی چندین مترجم ماشینی استفاده نمی‌شود، بلکه مرحله به مرحله ترجمه اولیه بهبود داده می‌شود. در ادامه به بررسی این روش و کارهای انجام‌شده در این زمینه می‌پردازیم.

در روش‌های متوالی ترجمه نهایی بر اساس ترجمه یک نوع مترجم ماشینی ایجاد می‌شود و نوع دیگر وظیفه کامل کردن این ترجمه و بهبود آن را دارد. در این روش، ترجمه نهایی می‌تواند بر پایه مترجم آماری و یا قاعده‌مند باشد. مقاله [7] با تکیه بر مترجم آماری، سعی کرده است که با ترجمه پیکره خود با استفاده از چندین مترجم قاعده‌مند، جدول عبارات^۱ مترجم آماری را کامل‌تر کند. هدف از این کار نیز قرارگرفتن ترجمه‌های قاعده‌مند در کنار ترجمه‌های مترجم آماری است که باعث کامل‌تر شدن جدول عبارات و در نهایت بهبود عملکرد مترجم آماری می‌شود. مقاله [4] نیز با تکیه بر مترجم آماری این کار را انجام می‌دهد. در این مقاله واژگان و عبارات جمله مبداء بر اساس زبان مقصد جابه‌جا می‌شوند، سپس با استفاده از جدول عباراتی که این جابه‌جایی در آن اعمال شده است، عمل ترجمه با استفاده از مترجم آماری یک‌نوا انجام می‌شود که باعث بهبود ترجمه نهایی از لحاظ ترتیب واژگان می‌شود. سانچز و همکاران [17] نیز جدول عبارات مترجم آماری را با استفاده از یک انتقال سطح پایین قاعده‌مند^۲ و فرهنگ لغت^۳ بهبود بخشیده‌اند. در این فرآیند هر عبارت در جدول عبارات با به‌کارگیری مجموعه‌ای از قوانین زبانی و نحوی تغییراتی در آن اعمال و سپس بررسی می‌شود که آیا عبارت جدید ترجمه‌ای در فرهنگ لغت مترجم قاعده‌مند دارد یا خیر، و در صورت وجود، ترجمه فرهنگ لغت را به جدول عبارات مترجم آماری اضافه می‌کند. تفاوت کارهای انجام‌شده در این زمینه با روش بیان‌شده در این مقاله، در انتخاب مترجم ماشینی پایه است. در کارهای معرفی‌شده از مترجم آماری به‌عنوان مترجم پایه استفاده و با استفاده از

اطلاعات مترجم قاعده‌مند ترجمه بهبود داده شده است؛ درحالی‌که در روش ارائه‌شده در این مقاله از مترجم قاعده‌مند به‌عنوان مترجم پایه استفاده شده است و با استفاده از اطلاعات مترجم آماری و رمزگشای یک‌نوا، ترجمه اولیه بهبود داده می‌شود.

مقالات [1] و [10] در کارهای خود فرهنگ لغت مترجم قاعده‌مند را با استفاده از جدول عبارات بهبود بخشیده‌اند. در مقالات نامبرده چون ترجمه‌های عبارات فرهنگ لغت مترجم قاعده‌مند به‌طور معمول ترجمه‌های روانی نیستند، سعی شده است با استفاده از روش‌هایی، معادل آن عبارات و یا خود آنها را در جدول عبارات پیدا کرده و ترجمه روان آن را به فرهنگ لغت اضافه کنند. تفاوت عملکرد این روش با روش ارائه‌شده در این مقاله در معماری و نحوه عملکرد رمزگشا است. در روش معرفی‌شده درواقع از مترجم قاعده‌مند استفاده شده که فرهنگ لغت آن بهبود یافته است. درحالی‌که در مقاله ارائه‌شده معماری و نحوه عملکرد رمزگشا به‌طور کامل متفاوت است که برتری آن نسبت به مقالات مطرح‌شده است، زیرا از یک رمزگشای یک‌نوا استفاده شده است که باعث روان‌تر شدن ترجمه نسبت به ترجمه‌های قاعده‌مند می‌شود. مقاله [2] نیز در کار خود به پس‌ویرایش^۴ خروجی مترجم قاعده‌مند پرداخته است. در این ایده یک پیکره موازی ارائه شده که یک سمت آن خروجی مترجم قاعده‌مند و سمت دیگر آن ترجمه بهبودیافته آن است. با استفاده از این پیکره یک جدول عبارات به‌وجود می‌آید و ترجمه خروجی مترجم قاعده‌مند به‌وسیله آن پس‌ویرایش و تغییراتی در آن ایجاد می‌شود که این کار باعث بهبود نتایج می‌شود. تفاوت این روش با روش ارائه‌شده در این مقاله نیز در معماری و نحوه عملکرد رمزگشا است. در این روش درواقع از مترجم ماشینی آماری با بهره‌مندی از یک جدول عبارات خاص و در روش ارائه‌شده در این مقاله از یک رمزگشای یک‌نوا با معماری متفاوت استفاده شده است.

معماری ترکیبی ارائه‌شده توسط اسپانا بونت و همکارانش در مقاله [8] به‌عنوان نزدیک‌ترین کار به فعالیت ما و شبیه به معماری ترکیبی ارائه‌شده در این مقاله است. مقاله یادشده ابتدا عمل ترجمه را با استفاده از مترجم قاعده‌مند انجام می‌دهد؛ سپس بر اساس درخت جمله ترجمه‌شده که از خروجی مترجم قاعده‌مند به‌دست آمده است، ترجمه عبارات مختلف جمله را تقویت می‌کنند. در مقاله معرفی‌شده ابتدا مرزهای عبارات مختلف را بر اساس خروجی مترجم قاعده‌مند به‌دست می‌آورند؛ سپس برای هر عبارت علاوه‌بر ترجمه

¹ Phrase table

² Shallow transfer rule

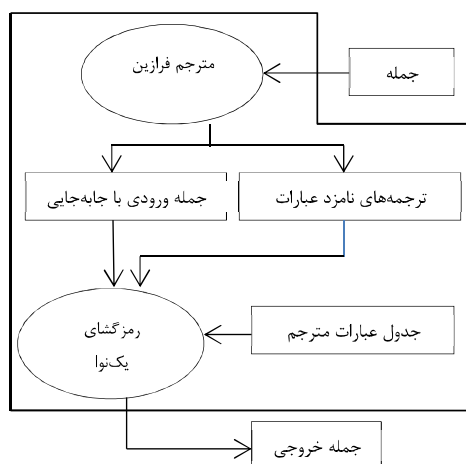
³ Dictionary

⁴ Post editing

هم‌ترازی بین جمله مبدا و مقصد در سطح عبارت است. علاوه بر این موارد، برای هر جمله مجموعه‌ای از اطلاعات دیگر مانند برجسب‌های اجزای کلام^۲، عبارات باهم‌آیی و پرکاربرد^۴ که نتیجه پردازش این مترجم است، نیز وجود دارد. مترجم آماری مبتنی بر عبارت استفاده‌شده در این معماری نیز، مترجم ماشینی موزز است که با استفاده از پیکره افک^۵ [9] آموزش داده شده و جدول عبارات آن استخراج شده است.

۱-۳- معماری ترکیبی

معماری ترکیبی ارائه‌شده بر پایه مترجم ماشینی فرازین است. شکل (۱) شمای کلی معماری را نشان می‌دهد. بر اساس این شکل، ابتدا جمله ورودی توسط مترجم فرازین ترجمه می‌شود و سپس خروجی آن مورد پردازش قرار می‌گیرد. در این پردازش ابتدا از بخش‌بندی عبارات ترجمه فرازین استفاده و مرز عبارات برای ترجمه مشخص می‌شود؛ سپس بر اساس ترتیب ترجمه عبارات، در ترجمه فرازین، این عبارات در کنار یکدیگر قرار می‌گیرند و ترتیب آنها با جمله مقصد یکسان می‌شود؛ در نتیجه می‌توان عمل ترجمه را با استفاده از یک رمزگشای یکنوا که در آن جابه‌جایی واژگان در نظر گرفته نمی‌شود، با هزینه ترجمه چندجمله‌ای انجام داد.



(شکل-۱): شمای کلی معماری ترکیبی
(Figure-1): hybrid architecture's schema

شکل (۲) یک نمونه از این تغییر ترتیب واژگان را نشان می‌دهد. هر یک از عبارات جمله مبدا بر اساس موقعیت قرارگیری ترجمه‌شان در ترجمه فرازین، در کنار یکدیگر قرار می‌گیرند.

مترجم قاعده‌مند، ترجمه‌های نامزدی از جدول عبارات انتخاب می‌شود؛ در نهایت از بین کلیه ترجمه‌های نامزد بر اساس مجموعه‌ای از ویژگی‌ها مانند مدل زبانی و امتیاز هر ترجمه، بهترین ترجمه انتخاب می‌شود. این فرآیند توسط رمزگشای مترجم ماشینی موزز^۱ [12] انجام می‌شود و مجموعه‌ای از ویژگی‌ها مانند احتمال ترجمه استخراج‌شده از فرهنگ لغت، از مترجم قاعده‌مند به آن اضافه شده است. این فرآیند بدون در نظر گرفتن مدل جابه‌جایی و به صورت یک‌نوا انجام می‌شود؛ ولی جابه‌جایی در ترتیب واژگان در جمله ورودی صورت نمی‌گیرد. همان‌طور که گفته شد، روش معرفی‌شده نزدیک‌ترین روش به روش ارائه‌شده در این مقاله است. تفاوت اصلی در جابه‌جایی عبارات جمله مبدا بر اساس ترجمه قاعده‌مند است. در مقاله [8] تغییری در ترتیب عبارات در جمله ورودی وجود ندارد و فقط از رمزگشای موزز با حذف مدل جابه‌جایی استفاده می‌شود؛ در حالی که در روش ارائه‌شده در این مقاله، ترتیب عبارات جمله ورودی بر اساس ترجمه قاعده‌مند تغییر می‌کند و سپس با استفاده از رمزگشای یک‌نوا عمل ترجمه انجام می‌شود که باعث می‌شود ترجمه نهایی از نظر ترتیب واژگان و عبارات در معیارهای مربوطه دارای عملکرد بهتری باشد.

کاری که ما در این مقاله انجام داده‌ایم، استفاده از یک رمزگشاهای یک‌نواخت برای عمل ترجمه و بهبود ترجمه مترجم قاعده‌مند با استفاده از جدول عبارات است. این رمزگشا در مقایسه با رمزگشای مترجم ماشینی موزز پیچیدگی زمانی کمتری دارد؛ همچنین برای احتمال‌دهی به ترجمه‌های مترجم قاعده‌مند خود که دارای احتمال ترجمه نیستند، روش‌هایی ارائه شده است که در بخش‌های بعد به تفصیل به آن می‌پردازیم.

۳- معماری سامانه

معماری ترکیبی ارائه‌شده بر اساس یک مترجم قاعده‌مند و یک مترجم آماری است. مترجم قاعده‌مند استفاده‌شده در این معماری، مترجم انگلیسی به فارسی فرازین^۲ است که توسط آزمایشگاه پردازش متن دانشگاه تهران تهیه و توسعه داده شده است. این مترجم ماشینی بر اساس قواعد نحوی و ساخت‌وازی، جمله ورودی را ترجمه می‌کند و خروجی آن علاوه بر ترجمه نهایی، حاوی اطلاعاتی شامل بخش‌بندی‌ها و ترتیب عبارات و

³ Port of speech tag

⁴ Simple phrase

⁵ AFEC

⁶ schema

¹ mooses

² <http://www.faraazin.ir/>

(شکل-۲): تغییر ترتیب کلمات جمله ورودی
(Figure-2): reordering of the input sentence

I go to home today .	جمله مبدأ
من امروز به خانه می‌روم.	ترجمه مقصد فرازین
I today to home go .	جمله مبدأ با جابه‌جایی بر اساس ترجمه مقصد

پایین قطر اصلی در فرآیند ترجمه، ترجمه نمی‌شوند و نحوه کامل کردن آن به صورت ستونی، از پایین به بالا است.

	I	today	to	home	go
I	↑	↑	↑	↑	↑
today		↑	↑	↑	↑
to			↑	↑	↑
home				↑	↑
go					↑

(شکل-۳): رمزگشای یک‌نوا بر پایه برنامه‌نویسی پویا
(Figure-3): monotone decoder based on dynamic programming

هر خانه $[i,j]$ بیان‌گر ترجمه عبارت i ام تا j ام تشخیص داده شده توسط فرازین است که احتمال آن بر اساس رابطه (۱) از ضرب وزن دار احتمال ترجمه و مدل زبانی محاسبه می‌شود. در این رابطه α و β فاکتورهای وزنی پارامترهای رابطه هستند که در بخش ۴ در رابطه با مقداردهی آنها صحبت می‌شود. در انتها نیز بهترین ده ترجمه به عنوان ترجمه عبارت مورد نظر در خانه $[i,j]$ ذخیره می‌شوند. انتخاب بهترین ده ترجمه برای هر خانه بر اساس روش‌های مکاشفه‌ای^۱ انجام شده است. دلیل انجام این کار این است که ترکیب ترجمه‌های نامزد غیر از بهترین ترجمه در هر خانه نیز بررسی شود تا بهترین ترکیب ترجمه‌ها به دست آید. نتایج آزمایش‌ها نشان می‌دهد که انتخاب بیش از ده ترجمه نامزد برای هر خانه تأثیری در نتایج ندارد.

$$p[i,j] = \max_{1 \leq k < j} [p[k+1,j] * \text{احتمال ترجمه } [i,k]] \quad (1)$$

([k+1,j] ترجمه + [i,k] ترجمه) [احتمال مدل زبانی] ^{β}

برای ترجمه هر خانه با اندیس^۲ $[i,j]$ عبارت مورد نظر به دو زیرعبارت^۳ در محل k شکسته می‌شود و ترجمه‌های خانه‌های $[i,k]$ و $[k+1,j]$ با هم ترکیب می‌شوند. مقدار k نیز برابر $j > k \geq i$ است. درواقع ترکیب‌های مختلف زیرعبارت، عبارت مورد نظر بررسی می‌شود تا ترکیب با بیشترین احتمال یافت شود. احتمال ترجمه این ترکیب نیز برابر با حاصل ضرب احتمالات ترجمه دو خانه ترکیب شده است و احتمال مدل زبانی نیز جداگانه محاسبه می‌شود. خانه‌های سطر نخست بیان‌گر ترجمه کل جمله از عبارت نخست تا عبارت مشخص شده هستند، و خانه $[0,n]$ حاوی ترجمه نهایی کل

به عنوان مثال، چون ترجمه عبارت “go” در انتهای ترجمه آمده است، پس در جابه‌جایی جمله مبدأ نیز، این عبارت به آخر جمله منتقل می‌شود. با ترجمه جمله جدید، درواقع ترتیب واژگان ترجمه نهایی بر اساس ترجمه فرازین صورت می‌گیرد. پس از جابه‌جایی جمله ورودی، فعالیت رمزگشای یک‌نوا آغاز می‌شود. در این مرحله، برای هر عبارت تشخیص داده شده توسط فرازین، از میان ترجمه‌های نامزد مترجم قاعده‌مند و آماری، بر اساس احتمال ترجمه و مدل زبانی بهترین ترجمه‌ها انتخاب می‌شود.

درواقع سامانه با استفاده از یک رمزگشای یک‌نوا که در بخش بعدی بیشتر آن را توضیح می‌دهیم، در هر مرحله، ترجمه‌های فرازین را با ترجمه‌های جدول عبارات ترکیب کرده و از بین ترجمه‌های نامزد، بهترین‌ها را بر اساس امتیاز هر ترجمه که در ادامه توضیح داده می‌شود، انتخاب می‌کند؛ و این کار را تا زمانی ادامه می‌دهد که در مرحله آخر بهترین ترجمه برای جمله مورد نظر بر اساس معیارهای معرفی شده در بخش ۲-۴ به دست آید.

۲-۳- رمزگشای یک‌نوا

ایده اصلی رمزگشای یک‌نوا برای نخستین بار در [18] ارائه شد. در این ایده برای جفت‌زبان‌هایی که از نظر ترتیب واژگان به هم نزدیک هستند، می‌توان با استفاده از برنامه‌نویسی پویا رمزگشایی طراحی کرد که عمل ترجمه را با هزینه چندجمله‌ای انجام دهد. معماری ترکیبی ما نیز بعد از ترجمه جمله توسط مترجم فرازین و در نخستین مرحله پردازش خود، جابه‌جایی جمله انگلیسی را با توجه به ترتیب واژگان ترجمه فرازین انجام می‌دهد. در نتیجه هم‌ترازی عبارات جمله مبدأ و مقصد یک‌نواخت می‌شود و می‌توان از رمزگشای یک‌نوا برای عمل ترجمه استفاده کرد. همان‌گونه که در شکل (۳) نیز قابل مشاهده است، برای عمل ترجمه از یک ماتریس دوبعدی با اندازه $n*n$ استفاده شده است که هر بعد آن بیان‌گر تعداد عبارات جمله مبدأ تشخیص داده شده توسط فرازین است. این ماتریس درواقع یک ماتریس بالامثلثی است که خانه‌های

¹ heuristic

² index

³ Sub-phrase

به عنوان واژگان مرزی ترجمه تا آن عبارت، ذخیره می شوند. درواقع می توان ماتریس ترجمه را یک ماتریس سه بعدی فرض کرد که بعد سوم آن مربوط به مدل زبانی است و اندازه آن برای هر ستون برابر تعداد واژگان مرزی ستون قبل در ماتریس ترجمه است. برای ترجمه هر خانه نیز، ترجمه را نسبت به تمامی واژگان مرزی ستون قبل بررسی می کنیم تا بهترین ترکیب با بیشترین احتمال یافت شود.

با توجه به شکل (۳)، به عنوان مثال برای شروع ابتدا خانه $[0,0]$ را ترجمه می کنیم. این خانه بیانگر ترجمه نخستین عبارت تشخیص داده شده توسط فرازین یعنی عبارت "i" به شرط نخست جمله است که با فرضیات بخش قبل، ترجمه می شود و بهترین ده ترجمه آن و دو کلمه آخر آنها به عنوان واژگان مرزی ستون نخست ذخیره می شوند؛ سپس خانه $[1,1]$ ترجمه می شود. در این خانه به ترجمه دومین عبارت تشخیص داده شده توسط فرازین، عبارت "today" به شرط دو کلمه ماقبل یعنی دو کلمه آخر ترجمه های نامزد عبارت "i" می پردازیم. در خانه $[0,1]$ نیز به ترکیب ترجمه عبارت نخست و دوم تشخیص داده شده توسط فرازین یعنی ترکیب ترجمه های خانه های $[0,0]$ و $[1,1]$ یعنی عبارت "i" و "today" پرداخته می شود و علاوه بر این، اگر کل عبارت خانه مورد نظر یعنی عبارت "i today" نیز در جدول عبارات وجود داشته باشد، از ترجمه آن استفاده می شود و این کار تا محاسبه بهترین ترجمه خانه $[0,n]$ ادامه پیدا می کند.

۳-۳- احتمال دهی و انتخاب ترجمه ها

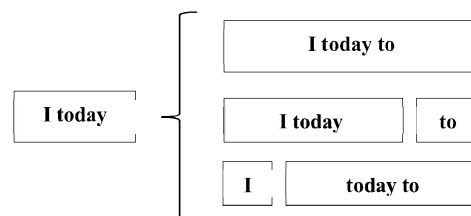
در فرآیند ترجمه هر عبارت، ابتدا در صورت وجود، ده ترجمه نخست آن عبارت از جدول عبارات استخراج می شود. بعد از آن نیز، در صورت وجود، ده ترجمه نخست فرازین، پس از محاسبه احتمال ترجمه برای آنها به ترجمه های استخراج شده از جدول عبارات اضافه می شوند. درنهایت از میان این ترجمه ها بر اساس احتمال ترجمه و مدل زبانی، بهترین ده ترجمه به عنوان ترجمه های آن عبارت انتخاب می شوند. ترجمه های فرازین دارای احتمال ترجمه نیستند و برای احتمال دهی به آنها از روش های مکاشفه ای زیر استفاده می کنیم:

- نخستین روشی که برای احتمال دهی به ترجمه های فرازین استفاده می شود، روش تطابق بین ترجمه فرازین و ترجمه های عبارت مورد نظر در جدول عبارات است. در صورتی که ترجمه فرازین در ترجمه های عبارت مورد نظر در جدول عبارات وجود داشته باشد از تطابق دقیق^۱ برای

^۱ Exactly matching

جمله است. درواقع در یک ستون زمانی که از قطر اصلی به سمت بالا پیش می رویم، به صورت عقب گرد از عبارت پایانی به سمت نخست جمله حرکت می کنیم.

خانه های غیر از قطر اصلی، نشانگر ترجمه ترکیب چندین عبارت هستند که برای ترجمه آنها از ترکیب ترجمه های زیرعبارات آن عبارت که درقبل محاسبه شده است، و جدول عبارات استفاده می شود. این فرآیند در شکل (۴) قابل مشاهده است. به عنوان مثال، برای ترجمه عبارت "i today to" که از سه زیرعبارت "i"، "today" و "to" تشکیل شده است، می توان کل عبارت را توسط جدول عبارات ترجمه کرد و یا می توان از ترکیب زیرعبارات "i" و "today to" با هم و یا ترکیب "i today" و "to" استفاده کرد که هر کدام از این خانه ها در مراحل قبل توسط بهترین ترجمه های آن عبارت تکمیل شده اند. می توان این فرآیند را به عنوان پیدا کردن بهترین مسیر از یک رأس به رأس دیگر در شبکه ترجمه در نظر گرفت. در بخش ۲ به شبکه ترجمه اشاره شد؛ فرآیند پیدا کردن بهترین مسیر در شبکه ترجمه به گونه ای است که اگر بین دو رأس یال مستقیمی وجود داشته باشد، برای پیمودن مسیر بین این رئوس، می توان از آن یال استفاده و همچنین می توان از رأس های بین این دو رأس برای پیمودن این مسیر نیز استفاده کرد.



(شکل-۴): حالات مختلف ترکیب زیرعبارات
(Figure-4): different way of combine sub-phrases

خانه های قطر اصلی این ماتریس نیز نشانگر ترجمه یک عبارت به تنهایی هستند؛ درنتیجه نمی توان برای ترجمه آنها از زیرعبارات استفاده کرد و فقط از ترجمه های فرازین و جدول عبارات استفاده می شود. احتمال خانه های قطر اصلی نیز بر اساس رابطه (۲) از ضرب وزن دار احتمال ترجمه و احتمال مدل زبانی به شرط دو کلمه ماقبل محاسبه می شود.

$$p[i, i] = \max_{\substack{\text{کلمات مرزی ستون قبل} \\ z \in}} [(i \text{ احتمال ترجمه})^\alpha * (z | i \text{ احتمال مدل زبانی})^\beta] \quad (2)$$

برای محاسبه احتمال مدل زبانی نیز، به ازای هر ستون، در تمامی خانه های آن، دو کلمه آخر بهترین ده ترجمه

۴- آزمایش‌ها

به‌منظور استفاده از جدول عبارات مترجم آماری، نیاز به یک پیکره موازی برای آموزش مترجم موزز است. پیکره استفاده‌شده برای این آموزش، پیکره دو زبانه افک است که اطلاعات این پیکره در جدول (۱) ارائه شده است.

از طرفی مجموع دادگان آزمون^۵ و توسعه^۶ ما نیز، یک مجموعه دادگان ۸۱۸ جمله‌ای با دامنه خبر است. به‌طور تقریبی نیمی از این مجموعه به‌منظور تنظیم پارامترهای^۷ مربوط به رابطه (۱ و ۲)، یعنی وزن‌های احتمال ترجمه و احتمال مدل زبانی، به‌عنوان دادگان توسعه و نیم دیگر آن به‌عنوان دادگان آزمون استفاده شده است. اطلاعات این مجموعه دادگان نیز در جدول (۱) قابل مشاهده است. این مجموعه دارای ۴ ترجمه مرجع^۸ است که برای محاسبه معیارهای ارزیابی از آنها استفاده شده است.

(جدول-۱): اطلاعات پیکره
(Table-1): corpus information

متوسط طول جملات (کلمه)	تعداد جملات	تعداد واژگان یکتا	تعداد واژگان	
۲۱	۶۸۳ هزار	۱۸۵ هزار	۵/۱۴ میلیون	دادگان آموزش انگلیسی
۲۲	۶۸۳ هزار	۲۰۲ هزار	۴/۱۵ میلیون	دادگان آموزش فارسی
۲۸	۴۱۸	۳/۰۵۱	۱۱/۶۲۶	دادگان توسعه
۲۷	۴۰۰	۲/۷۶۴	۱۰/۶۵۸	دادگان آزمون

۴-۱- معرفی سامانه‌های ترکیبی

دو پارامتر تعیین‌کننده در کیفیت خروجی یک مترجم ماشینی، ارائه ترتیب مناسب برای واژگان و انتخاب معادل‌های مناسب برای واژگان و عبارات است. به‌منظور مقایسه عملکرد سامانه‌های مترجم ماشینی مختلف، سامانه‌هایی با ترکیب‌های مختلف این دو پارامتر پیاده‌سازی شده‌اند که اطلاعات آنها در جدول (۲) قابل مشاهده است. به‌منظور مقایسه عملکرد سامانه ترکیبی نسبت به دو سامانه قاعده‌مند و آماری، سامانه یک

احتمال‌دهی به ترجمه فرازین استفاده می‌شود و احتمال آن به ترجمه فرازین اختصاص می‌یابد. در صورت عدم تطابق دقیق از روش کم‌ترین فاصله ویرایش^۱ برای تطابق جزئی^۲ استفاده می‌شود. در این روش نیز ترجمه‌ای که دارای کمترین فاصله ویرایش است، احتمالش به ترجمه فرازین اختصاص داده می‌شود.

• در صورتی که ترجمه‌ای از فرازین مطابقتی در جدول عبارات پیدا نکند از تعدادی فرضیات مکاشفه‌ای برای احتمال‌دهی به ترجمه‌های فرازین استفاده می‌شود. به‌عنوان مثال برای احتمال‌دادن به ترجمه دوم فرازین از احتمال ترجمه دوم، ترجمه‌های آماری استخراج‌شده از جدول عبارات برای عبارت مورد نظر استفاده می‌شود و یا مجموع احتمالات ترجمه‌های آماری استخراج‌شده برای عبارت مورد نظر به‌صورت یکسان بین ترجمه‌های فرازین تقسیم می‌شود.

همان‌طور که در قبل اشاره شد، برای هر عبارت، بهترین ده ترجمه از بین ترجمه‌های نامزد آن عبارت انتخاب و برای این منظور، از روش رأی‌گیری^۳ استفاده می‌شود. در این روش که به‌طور معمول در مترجم‌های ماشینی ترکیبی با معماری موازی استفاده می‌شود، ترجمه‌هایی که در چندین مترجم ماشینی، به‌صورت یکسان وجود دارند، یک ترجمه در نظر گرفته می‌شوند و احتمال ترجمه همگی آنها با هم جمع می‌شود. به‌عنوان مثال در شبکه ترجمه که پیش‌تر به آن اشاره کردیم، یال‌هایی که دارای ترجمه یکسان هستند، به یک یال با احتمال ترجمه‌ای برابر با جمع احتمالات ترجمه آنها تبدیل می‌شوند. هدف از این کار نیز این است که اگر یک ترجمه، توسط چندین مترجم ماشینی تولید شده باشد، به‌احتمال ترجمه مناسبی است و باید احتمال بالاتری نسبت به بقیه ترجمه‌های نامزد به آن اختصاص داده شود؛ از این‌رو احتمالات ترجمه‌ای که مترجم‌های ماشینی مختلف به آن ترجمه اختصاص داده‌اند، با هم جمع می‌شوند. گفتنی است در این حالت امکان دارد احتمال ترجمه به امتیاز ترجمه تبدیل شود. به این منظور ما نیز زمانی که ترجمه‌های جدول عبارات استخراج شدند و ترجمه‌های فرازین نیز احتمال ترجمه‌شان محاسبه شد، این عمل را انجام داده و سپس بر اساس احتمال ترجمه‌های موجود و احتمال مدل زبانی آنها نسبت به دو کلمه قبل از خود با استفاده از یک مدل زبانی^۳-گرام^۴، بهترین ده ترجمه را به‌عنوان ترجمه آن عبارت در نظر می‌گیریم و ذخیره می‌کنیم.

^۵ Test data

^۶ development data

^۷ Parameter tuning

^۸ Gold data

^۱ Minimum edit distance

^۲ Partial matching

^۳ voting

^۴ trigram

به عنوان سامانه قاعده مند و سامانه دو به عنوان سامانه آماری پیاده سازی شده است.

(جدول ۲): اطلاعات سامانه ها

(Table-2): systems information

نام سامانه	ترتیب کلمات	انتخاب لغات
۱	مترجم قاعده مند	مترجم قاعده مند
۲	مترجم آماری	مترجم آماری
۳	مترجم قاعده مند	مترجم آماری
۴	مترجم قاعده مند	مترجم قاعده مند و آماری
۵	مترجم آماری	مترجم قاعده مند و آماری

بر اساس این جدول، سامانه یک، یک سامانه قاعده مند بوده که در واقع مترجم فرازین است. مترجم فرازین به منظور انتخاب لغات از واژگان داخلی خود و از روشی مبتنی بر قواعد نحوی برای ترتیب واژگان ترجمه استفاده می کند. در پیاده سازی سامانه دو که یک سامانه آماری است، فقط از اطلاعات مترجم آماری مبتنی بر عبارت موز استفاده شده است. این مترجم برای انتخاب لغات از جدول عبارتی که به صورت آماری از پیکره موازی استخراج کرده است و برای ارائه ترتیب و ترجمه از مدل جابه جایی لغوی و مبتنی بر فاصله استفاده می کند. در این پیاده سازی ابتدا دادگان توسط مترجم موز ترجمه شده است و سپس از هم ترازی در سطح عبارت ترجمه آن استفاده شده تا ترتیب واژگان ترجمه و رمز عبارات تشخیص داده شود. همچنین برای مرحله انتخاب لغات فقط از جدول عبارات این مترجم استفاده شده است.

به منظور مقایسه ترتیب واژگان در ترجمه های مترجم قاعده مند و آماری، سامانه ترکیبی، با بهره مندی از ترتیب واژگان مترجم قاعده مند در قالب سامانه چهار و با بهره مندی از ترتیب واژگان مترجم آماری در قالب سامانه پنج پیاده سازی شده، که در هر پیاده سازی از هم ترازی در سطح عبارت خروجی مترجم مربوطه، یعنی مترجم فرازین و موز استفاده شده است.

برای مقایسه اینکه برای انتخاب لغات بهتر است از یک منبع استفاده شود یا هر دو، سامانه با ترتیب واژگان مترجم قاعده مند یکبار در حالتی پیاده سازی شده که برای انتخاب لغات فقط از اطلاعات مترجم آماری یعنی جدول عبارات مترجم موز استفاده شده است که به عنوان سامانه سه معرفی

شده، و بار دیگر در حالت ترکیبی، که همان سامانه چهار است، و از هر دو منبع مترجم قاعده مند و آماری، یعنی جدول عبارات مترجم موز و ترجمه های فرازین استفاده می کند. با پیاده سازی این سامانه ها می توان سامانه قاعده مند را که فقط از ترجمه های مترجم قاعده مند برای انتخاب لغات استفاده می کند، با سامانه ای که برای این عمل فقط از ترجمه های مترجم آماری استفاده می کند و همچنین سامانه ترکیبی مقایسه کرد.

۲-۴- ارزیابی، نتایج و تحلیل

برای ارزیابی نتایج سامانه ها از سه معیار BLEU [15]، NIST [6] و LR-Score [3]، استفاده شده است. معیار بلو به عنوان رایج ترین معیار ارزیابی نتایج مترجم های ماشینی است که بر پایه تطابق چند-گرام های^۱ متفاوت در ترجمه سامانه و ترجمه مرجع است. همچنین معیار NIST نیز بر پایه معیار بلو است، البته با تغییراتی در وزن دهی به چندگرام های متفاوت برای امتیازدهی نهایی. معیار LR-Score نیز با استفاده از هم ترازی ترجمه سامانه و ترجمه مرجع، موقعیت ترجمه هر عبارت در ترجمه سامانه را نسبت به ترجمه مرجع بررسی کرده و در نهایت با مقایسه موقعیت کلیه عبارات و اختلاف بین این موقعیت ها بین ترجمه سامانه و ترجمه مرجع، ترتیب واژگان ترجمه سامانه نسبت به ترتیب واژگان ترجمه مرجع را می سنجد.

به منظور احتمال دهی به ترجمه های فرازین در روش تطابق جزئی در بخش ۳-۳، برای ترجمه هایی که طول آنها کمتر از ۴ است، اختلاف یک و برای طول بیشتر از چهار، اختلاف ۴۰٪ طول ترجمه به عنوان فاصله ویرایش در روش کم ترین فاصله ویرایش، در نظر گرفته می شود. البته بررسی نتایج آزمایش ها نشان می دهد که تغییر میزان تطابق، تأثیر زیادی در نتایج سامانه ندارد؛ همچنین برای احتمال دهی به ترجمه نخست فرازین، اگر این ترجمه برابر با خود عبارت باشد، یعنی فرازین نتوانسته باشد عبارت مورد نظر را ترجمه کند، کمترین احتمال ترجمه های آماری و در غیر این صورت بیشترین احتمال ترجمه های آماری به آن اختصاص داده می شود. اگر عبارتی در جدول عبارات وجود نداشته باشد نیز مجموع احتمالات را برابر ۰/۵، بیشترین احتمال را برابر ۰/۳ و کمترین احتمال را ۰/۲ در نظر می گیریم.

به منظور تنظیم کردن پارامترهای رابطه (۱)، سامانه های موجود را با مقادیر مختلف پارامترهای α و β بر

^۱ N-gram

روی دادگان توسعه اجرا می‌کنیم. پارامترهای سامانه با بیشترین مقدار معیار بلو به‌عنوان پارامترهای اصلی هر سامانه انتخاب می‌شوند. وزن پارامترهای هر سامانه در جدول (۳) قابل مشاهده است.

(جدول-۳): وزن پارامترهای سامانه

(Table-3): weight of system parameters

نام سامانه	۲	۳	۴	۵
وزن احتمال ترجمه α	0.62	0.48	0.59	0.58
وزن احتمال مدل زبانی β	0.38	0.52	0.41	0.42

پس از تنظیم پارامترهای هر سامانه، آنها را بر روی دادگان آزمون اجرا کردیم. نتایج این آزمایش‌ها در جدول (۴)

(جدول-۴): نتایج آزمایش‌ها

(Table-4): Result of the experiments

نام سامانه	۱	۲	۳	۴	۵
توضیحات سامانه	سامانه قاعده‌مند	سامانه آماری	سامانه قاعده‌مند با انتخاب لغات مترجم آماری	سامانه ترکیبی با ترتیب واژگان مترجم قاعده‌مند	سامانه ترکیبی با ترتیب واژگان مترجم آماری
BLEU	18.66	22.50	12.39	22.36	22.14
NIST	7.44	7.79	6.32	8.13	7.77
LR-Score	0.52	0.46	0.52	0.50	0.46

سامانه‌ای بهتر از سامانه آماری، فرضیاتی را به سامانه ترکیبی اضافه می‌کنیم.

در این فرضیات با توجه به نقاط قوت عملکرد مترجم فرازین و موزز، بعضی عبارات فقط توسط یک مترجم، ترجمه می‌شود. از جمله عباراتی که فقط توسط مترجم فرازین ترجمه می‌شوند، عبارات باهم‌آیی و پرکاربرد و عبارات خاص هستند که با برچسب اجزای کلام اسم خاص^۱ شناخته می‌شوند و عبارات باهم‌آیی و پرکاربرد نیز به‌طور معمول از ترکیب چندین کلمه به‌وجود می‌آیند و به عباراتی اطلاق می‌شوند که فرازین برای آنها ترجمه مشخص و یک‌پارچه‌ای دارد، مانند واژگان خاص چندبخشی. در آزمایش‌هایی که انجام گرفت، مشخص شد که اگر برای ترجمه این عبارات فقط از ترجمه‌های مترجم فرازین استفاده شود، نتایج سامانه ترکیبی بهبود می‌یابد؛ همچنین برای بررسی نقاط قوت مترجم موزز، آزمایش‌هایی در زمینه برچسب‌های اجزای کلام انجام شد. در ابتدا خروجی مترجم فرازین، موزز و سامانه ترکیبی با استفاده از یک برچسب‌گذار اجزای کلام^۲ فارسی، برچسب زده شدند. در

برای مقایسه ترتیب واژگان مترجم قاعده‌مند و آماری نیز، همان‌گونه که مشخص است در تمامی معیارها، به‌خصوص LR-Score نتایج ترتیب واژگان سامانه آماری ضعیف‌تر از قاعده‌مند است. همچنین در این معیار نتایج سامانه قاعده‌مند بیشترین مقدار را دارد و سامانه آماری دارای کمترین مقدار است. این موضوع نشان‌دهنده آن است که اگر برای انتخاب پارامتر ترتیب واژگان از مترجم قاعده‌مند استفاده شود، نتایج، بهتر از حالتی است که از مترجم آماری استفاده شود. انتخاب مترجم قاعده‌مند به‌عنوان سامانه پایه در معماری ترکیبی ارائه‌شده، انتخاب صحیحی است. همچنین سامانه ترکیبی که بر پایه ترتیب واژگان مترجم قاعده‌مند است، در این معیار دارای عملکرد بهتری نسبت به سامانه ترکیبی است که بر پایه ترتیب کلمات مترجم آماری است.

نتایج آزمایش‌ها نشان می‌دهد که ترتیب واژگان مترجم قاعده‌مند بهتر از مترجم آماری است و اگر برای انتخاب لغات از هر دو منبع استفاده شود، بهترین سامانه ترکیبی به‌دست می‌آید. به‌منظور ارتقای سامانه و دستیابی به

¹ Proper noun

² POS tagger

با اعمال این فرضیات بر روی سامانه ترکیبی، نتایج بدست آمده و مقایسه آن با رقیب اصلی یعنی سامانه آماری در جدول (۵) مشاهده می‌شود. همان‌گونه که مشخص است در تمامی معیارها سامانه ترکیبی بهتر از سامانه آماری و در بیش‌تر آنها از سامانه قاعده‌مند نیز بهتر است. همچنین آزمون معناداری^۲ [11] بر روی نتایج این دو سامانه انجام شد و نتایج نشان داد که تفاوت این دو سامانه معنادار بوده است و وابسته به پیکره و مجموعه دادگان نیست.

(جدول-۵): نتایج نهایی و آزمایش معناداری
(Table-5): final results and significance test

مقدار p-value آزمایش معناداری	معیار BLEU	سامانه
0.05	23.57	سامانه ترکیبی با ترتیب کلمات مترجم قاعده‌مند
	22.50	سامانه آماری

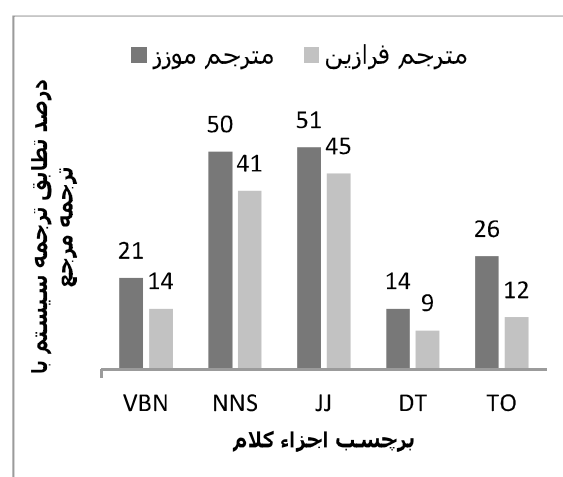
(جدول-۶): نمونه‌ای از نتایج
(Table-6): result sample

معیار بلو برای یک جمله	ترجمه	سامانه‌ها
جمله انگلیسی	four afghans, including two students, were also killed, said hashmat stanikzai, spokesman for kabul's police chief.	جمله انگلیسی
ترجمه مرجع	چهار افغان از جمله دو دانش‌آموز نیز طبق گفته‌های حشمت استانی‌کزی، سخنگوی رییس پلیس کابل کشته شده‌اند.	ترجمه مرجع
ترجمه سامانه قاعده‌مند	، حشمت stanikzai، سخنگو برای رییس پلیس کابل گفتند افغان‌ها چهار، از جمله دو تن از دانشجویانی، همچنین کشته شدند.	0.11
ترجمه سامانه آماری	چهار افغان‌ها از جمله دو دانشجو کشته نیز گفت حشمت استانی‌کزی سخنگوی رییس پلیس کابل.	0.21
ترجمه سامانه ترکیبی	، حشمت استانی‌کزی، سخنگوی رییس پلیس کابل گفت: چهار افغان از جمله دو دانشجو، همچنین کشته شدند.	0.36

در جدول (۶) نیز یک نمونه از ترجمه بهبودیافته توسط سامانه ترکیبی در مقایسه با سامانه آماری و قاعده‌مند با

مرحله بعد برای هر برچسب^۱، بررسی شد که چند درصد از ترجمه‌های آن برچسب در هر سامانه، تطابق بیشتری با ترجمه‌های مرجع دارد. بررسی آزمایش‌ها نشان داد که مترجم موز در بیش‌تر برچسب‌ها تطابق بیشتری با ترجمه‌های مرجع آن برچسب دارد. نتایج این آزمایش‌ها نشان‌دهنده این است که عملکرد مترجم موز در ترجمه این دسته از برچسب‌های اجزای کلام، بهتر از عملکرد مترجم فرازین است و اگر برای ترجمه آنها فقط از ترجمه‌های مترجم موز استفاده شود، می‌توان عملکرد سامانه ترکیبی را بهبود داد. از میان این دسته برچسب‌های اجزای کلام، برچسب‌هایی که دارای فراوانی کمی بودند، به‌عنوان مثال تعداد آن برچسب در کل پیکره برابر ده عدد بود، حذف شده و برچسب‌هایی که دارای فراوانی بالایی بودند و علاوه‌براین عملکرد سامانه ترکیبی نیز در ترجمه آنها و تطابق با ترجمه مرجع پایین بود، انتخاب شدند. برای هر برچسب در این دسته درحالتی که برای ترجمه آنها فقط از ترجمه‌های مترجم موز استفاده شود و حالت معمولی که همان حالت ترکیبی است، عملکرد سامانه ترکیبی بررسی شد. برچسب‌هایی که باعث بهبود عملکرد سامانه ترکیبی شدند به‌همراه میزان تطابق ترجمه آنها در هر سامانه با ترجمه مرجع در شکل (۵) قابل مشاهده است.

در نتیجه برای ترجمه این دسته از برچسب‌های اجزای کلام، فقط از ترجمه‌های جدول عبارات مترجم موز استفاده شد. این برچسب‌ها شامل VBN برچسب مربوط به افعال در زمان گذشته، NNS برچسب اسامی جمع، JJ برچسب مربوط به صفات، DT برچسب مربوط به حروف اشاره و TO برچسب مربوط به حرف اضافه to است.



(شکل-۵): نتایج آزمایش‌های برچسب اجزاء کلام
(Figure-5): results of POS tag experiment

² Significance test

¹ tag

مناسب نیز می‌توان علاوه بر احتمال ترجمه و مدل زبانی، ویژگی‌های دیگری مانند فرضیات اضافه شده به سامانه ترکیبی، برای رمزگشا تعریف کرد تا مرحله انتخاب ترجمه مناسب نیز هوشمندانه‌تر شود.

6- References

۶- مراجع

- [1] A. Antonova and A. Misyurev, "Improving the precision of automatically constructed human-oriented translation dictionaries", in *Proceedings of the 3rd Workshop on Hybrid Approaches to Translation (HyTra)@ EACL*, 2014, pp. 58–66.
- [2] H. Béchara, R. Rubino, Y. He, Y. Ma, and J. van Genabith, "An Evaluation of Statistical Post-Editing Systems Applied to RBMT and SMT Systems,," in *COLING*, 2012, vol. 21, pp. 5–230.
- [3] A. Birch and M. Osborne, "LRscore for evaluating lexical and reordering quality in MT,," in *Proceedings of the Joint Fifth Workshop on Statistical Machine Translation and Metrics MATR*, 2010, pp. 327–332.
- [4] Y. Chen and A. Eisele, "Hierarchical hybrid translation between english and german,," in *Proceedings of the 14th Annual Conference of the European Association for Machine Translation*, 2010, pp. 90–97.
- [5] M. R. Costa-Jussa, M. Farrús, J. B. Marino, and J. A. R. Fonollosa, "Study and Comparison of Rule-Based and Statistical Catalan-Spanish Machine Translation Systems,," *Comput. informatics*, vol. 31, no. 2, pp. 245–270, 2012.
- [6] G. Doddington, "Automatic evaluation of machine translation quality using n-gram co-occurrence statistics,," in *Proceedings of the second international conference on Human Language Technology Research*, 2002, pp. 138–145.
- [7] A. Eisele, C. Federmann, H. Uszkoreit, H. Saint-Amand, M. Kay, M. Jellinghaus, S. Hunsicker, T. Herrmann, and Y. Chen, "Hybrid Architectures for Multi-Engine Machine Translation,," *Proc. Transl. Comput. 30. ASLIB Transl. Comput. Conf. (TATC)*, 30th, Novemb. 27-28, London, United Kingdom, 2008.
- [8] C. España Bonet, L. Márquez Villodre, G. Labaka, A. de Ilaraza Sánchez, and K. Sarasola Gabiola, "Hybrid machine translation guided by a rule-based system,," in *Machine translation summit XIII: proceedings of the 13th machine translation summit*, September 19-23, 2011, Xiamen, China, 2011, pp. 554–561.
- [9] S. FattanehJabbari and S. M. M. Ziabary, "Developing an open-domain English-Farsi translation system using AFEC: Amirkabir bilingual Farsi-English corpus,," in *The Fourth*

معیاری مشابه معیار بلو برای یک جمله^۱ را مشاهده می‌کنیم. همان‌گونه که مشخص است، سامانه آماری به‌دلیل اینکه تطابق بیشتری با ترجمه مرجع در چند-گرام‌های متفاوت دارد، عملکرد بهتری در معیار مربوطه نیز دارد، ولی همان‌گونه که مشخص است، عملکرد سامانه قاعده‌مند در زمینه ترتیب واژگان و ترجمه مناسب افعال، بهتر از عملکرد سامانه آماری است و سامانه ترکیبی که از نقاط قوت هر دو سامانه بهره برده است، دارای بهترین ترجمه است. در این ترجمه واژگان پررنگ از سامانه قاعده‌مند و بقیه واژگان و عبارات از سامانه آماری انتخاب شده است. همچنین ترتیب واژگان و عبارات نیز برگرفته از سامانه قاعده‌مند است.

۵- نتیجه‌گیری

در این مقاله با استفاده از یک رمزگشای یکنوا با هزینه چند جمله‌ای، یک مترجم ماشینی ترکیبی طراحی کردیم. این مترجم بر پایه یک مترجم قاعده‌مند است که ترتیب واژگان در ترجمه آن بر اساس مترجم قاعده‌مند و انتخاب لغات آن به‌صورت ترکیبی از مترجم‌های قاعده‌مند و آماری است. آنچه که از نتایج آزمایش‌ها مشخص است، این است که مواردی در ترجمه مانند ترتیب واژگان که مربوط به قواعد نحوی زبان است، بهتر است، توسط مترجم‌های ماشینی قاعده‌مند پوشش داده شوند. در معیارهایی که موارد نحوی مانند ترتیب واژگان را بررسی می‌کنند، این مترجم‌ها عملکرد بهتری نسبت به مترجم‌های آماری دارند؛ ولی در مقابل این مترجم‌ها در به‌کارگیری لغات مناسب و نزدیک به ترجمه انسانی ضعیف هستند و در معیارهایی مانند بلو عملکرد بدتری نسبت به مترجم‌های آماری دارند؛ درنتیجه برای رسیدن به سامانه بهینه باید برای انتخاب لغات، از ترجمه‌های نامزد مترجم قاعده‌مند و آماری با هم استفاده کرد؛ پس می‌توان با ترکیب این دو نوع معماری به مترجم ماشینی رسید که در همه معیارها بتواند فاصله اختلاف این دو نوع مترجم ماشینی را کاهش دهد و از نقاط قوت هر دو نوع بهره ببرد.

معماری ترکیبی را می‌توان با استفاده از رمزگشاهای دیگری مانند رمزگشای مترجم ماشینی موزز با ویژگی‌های بیشتر و کامل‌تر، پیاده‌سازی کرد. همچنین می‌توان با افزایش تعداد مترجم‌های ماشینی و انتخاب توابع هوشمندتر برای احتمال‌دهی به ترجمه‌های مترجم‌های قاعده‌مند و انتخاب ترجمه مناسب از بین ترجمه‌های نامزد موجود عملکرد معماری‌های ترکیبی را افزایش داد. برای انتخاب ترجمه

^۱ Sentence level BLEU



حکیمه فدایی

دانشجوی مقطع دکترای رشته مهندسی نرم افزار در دانشگاه تهران است. همچنین ایشان مدرک کارشناسی و کارشناسی ارشد خود را در رشته مهندسی نرم افزار در دانشگاه شهید بهشتی کسب کرده اند و زمینه پژوهشی ایشان پردازش زبان طبیعی و به طور خاص مترجم ماشینی است.

نشانی رایانامه ایشان عبارت است از:

h.fadaei@ut.ac.ir



هشام فیلی

تحصیلات خود را در مقطع کارشناسی مهندسی نرم افزار در دانشکده مهندسی کامپیوتر دانشگاه صنعتی شریف با رتبه یک در سال ۱۳۷۶ به پایان رساند؛

سپس مقاطع کارشناسی ارشد نرم افزار و دکترای هوش مصنوعی را به ترتیب در سال های ۱۳۷۸ و ۱۳۸۵ در همان دانشکده تکمیل کرد. از سال ۱۳۸۷ تاکنون عضو هیأت علمی دانشکده مهندسی برق و کامپیوتر دانشکده فنی دانشگاه تهران است. زمینه های پژوهشی مورد علاقه ایشان پردازش هوشمند متن و زبان طبیعی، ترجمه ماشینی، داده کاوی، بازیابی اطلاعات و شبکه های اجتماعی هستند.

نشانی رایانامه ایشان عبارت است از:

hfaily@ut.ac.ir

Workshop on Computational Approaches to Arabic Script-based Languages, 2012, pp. 17.

- [10] N. Habash, B. Dorr, and C. Monz, "Symbolic-to-statistical hybridization: extending generation-heavy machine translation," *Mach. Transl.*, vol. 23, no. 1, pp. 23–63, 2009.
- [11] P. Koehn, "Statistical Significance Tests for Machine Translation Evaluation," in *EMNLP*, 2004, pp. 388–395.
- [12] P. Koehn, H. Hoang, A. Birch, C. Callison-Burch, M. Federico, N. Bertoldi, B. Cowan, W. Shen, C. Moran, R. Zens, and others, "Moses: Open source toolkit for statistical machine translation," in *Proceedings of the 45th annual meeting of the ACL on interactive poster and demonstration sessions*, 2007, pp. 177–180.
- [13] G. Labaka, C. España-Bonet, L. Márquez, and K. Sarasola, "A hybrid machine translation architecture guided by syntax," *Mach. Transl.*, vol. 28, no. 2, pp. 91–125, 2014.
- [14] W. Macherey and F. J. Och, "Consensus translations from multiple machine translation systems." Google Patents, 2012.
- [15] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, "BLEU: a method for automatic evaluation of machine translation," in *Proceedings of the 40th annual meeting on association for computational linguistics*, 2002, pp. 311–318.
- [16] E. Park, O. Kwon, K. Kim, and Y. Kim, "A Classification-based Approach for Hybridizing Statistical Machine Translation and Rule-based Machine Translation," pp. 1–9, 2015.
- [17] V. M. Sánchez-Cartagena, F. Sánchez-Martínez, J. A. Pérez-Ortiz, and others, "Integrating shallow-transfer rules into phrase-based statistical machine translation," 2011.
- [18] C. Tillmann, S. Vogel, H. Ney, and A. Zubiaga, "A DP based search using monotone alignments in statistical translation," in *Proceedings of the eighth conference on European chapter of the Association for Computational Linguistics*, 1997, pp. 289–296.

حسین خاتمی

دانشجوی کارشناسی ارشد رشته مهندسی نرم افزار در دانشگاه تهران از سال ۱۳۹۳ است. همچنین مقطع کارشناسی را از دانشگاه صنعتی سجاد مشهد در سال ۱۳۹۲ در رشته مهندسی



نرم افزار اخذ کرده است. زمینه پژوهشی ایشان پردازش زبان طبیعی و به طور خاص مترجم ماشینی است.

نشانی رایانامه ایشان عبارت است از:

h.khatami@ut.ac.ir

