

مجموعه مقالات چهارمین همایش ملی زبان‌شناسی رایانشی

---

به کوشش  
مسعود قیومی  
آذین شهریاری‌فرد



|                                |  |
|--------------------------------|--|
| سرشناسه<br>عنوان و نام پدیدآور | همایش ملی زبان‌شناسی رایانشی (چهارمین : ۱۳۹۶ : تهران)<br>مجموعه مقالات چهارمین همایش ملی زبان‌شناسی رایانشی/بسه کوشش مسعود قیسومی،<br>آذین شهریاری‌فرد؛ برگزارکننده انجمن زبان‌شناسی ایران؛ با همکاری پژوهشگاه علوم<br>انسانی و مطالعات فرهنگی، پایگاه استنادی علوم جهان اسلام (ISC)، نشر نویسه پاریسی.<br>تهران: نشر نویسه پاریسی، ۱۳۹۶.<br>۲۰۹ ص: مصور، جدول، نمودار.<br>۹۷۸-۶۰۰-۷۰۳۰-۵۹-۲ |
| مشخصات نشر                     | تهران: نشر نویسه پاریسی، ۱۳۹۶.   |
| مشخصات ظاهری                   | ۲۰۹ ص: مصور، جدول، نمودار.   |
| شابک                           | ۹۷۸-۶۰۰-۷۰۳۰-۵۹-۲  |
| وضعیت فهرست نویسی              | فیبا   |
| موضوع                          | زبان‌شناسی کامپیوتری -- کنگره‌ها   |
| موضوع                          | Computational (Linguistics) -- Congresses  |
| موضوع                          | زبان‌شناسی -- ایران -- کنگره‌ها  |
| موضوع                          | Linguistics -- Iran -- Congresses  |
| شناسه افزوده                   | قیسومی، مسعود، ۱۳۵۸ - ، گردآورنده.   |
| شناسه افزوده                   | شهریاری‌فرد، آذین، ۱۳۶۶ - ، گردآورنده.   |
| شناسه افزوده                   | منظوری، ریحانه، ۱۳۶۳ - ، ویراستار.   |
| شناسه افزوده                   | انجمن زبان‌شناسی ایران.  |
| رده بندی کنگره                 | ۱۳۹۶ ۸۸۸ت / PIR۲۶۵۳  |
| رده بندی دیویی                 | ۴۶۰/۲۸   |
| شماره کتابشناسی ملی            | ۵۰۶۶۰۶۶  |



تهران، صندوق پستی ۱۳۷۹-۱۶۷۶۵  
 تلفن: ۷۷۰۵۳۲۴۶  
 شماره: ۷۷۰۵۳۲۴۶  
 سامانه پیام کوتاه: ۳۰۰۰۴۵۵۴۵۵۴۱۴۲  
 وبگاه نشر نویسه پاریس:  
[www.neveeseh.com](http://www.neveeseh.com)



انجمن زبان‌شناسی ایران

تهران، بزرگراه چمران، پل مدیریت، خیابان علامه  
 طباطبایی جنوبی، دانشکده ادبیات فارسی و  
 زبان‌های خارجی دانشگاه علامه طباطبایی. طبقه  
 اول، اتاق ۱۱۷. صندوق پستی: ۱۵۹۷۶۳۳۱۱۱  
 تلفن: ۸۸۶۹۰۰۲۲ شماره: ۸۸۶۹۰۰۲۲  
[www.lsi.ir](http://www.lsi.ir)

### دارای درجه علمی - پژوهشی

#### نمایه شده در پایگاه استنادی علوم جهان اسلام (ISC)

دارای درجه علمی - پژوهشی، بر اساس بند ۱ ماده ۳ آیین‌نامه نحوه برگزاری و ساماندهی همایش‌های علمی ابلاغیه شماره ۳/۱۶۸۲۱۶ مورخ ۱۳۹۳/۹/۱۱ وزارت علوم، تحقیقات و فن‌آوری.

همه حقوق محفوظ و متعلق به «نشر نویسه پاریس» است.

تکثیر، انتشار و ترجمه این اثر یا قسمتی از آن به هر شیوه، بدون مجوز قبلی و کتبی ممنوع و مورد پیگرد قانونی قرار خواهد گرفت.

شابک: ۹۷۸-۶۰۰-۷۰۳۰-۵۹-۲

ISBN: 978-600-7030-59-2

#### مجموعه مقالات چهارمین همایش ملی زبان‌شناسی رایانشی

|                             |  |
|-----------------------------|--|
| به کوشش                     | مسعود قیومی - آذین شهرباری فرد                                     |
| ویراستار چکیده‌های انگلیسی  | ریحانه منظوری  |
| طرح جلد و یونیفورم          | محمد محرابی <a href="http://www.mehraabi.com">www.mehraabi.com</a> |
| صفحه‌آرایی و آماده‌سازی چاپ | محمد محرابی <a href="http://www.mehraabi.com">www.mehraabi.com</a> |
| چاپ و صحافی                 | روز  |
| شمارگان                     | ۲۰۰ نسخه   |
| نوبت چاپ                    | اول، ۱۳۹۶  |
| قیمت                        | ۲۰۰۰ تومان   |

# مجموعه مقالات چهارمین همایش ملی زبان‌شناسی رایانشی

## برگزارکننده

انجمن زبان‌شناسی ایران

## با همکاری

پژوهشگاه علوم انسانی و مطالعات فرهنگی

پایگاه استنادی علوم جهان اسلام (ISC)

نشر نویسه پارسی

## سازمان همایش

دکتر مسعود قیومی دبیر علمی

آذین شهریاری فرد دبیر اجرایی

## کمیته علمی

دکتر سیدمصطفی عاصی عضو هیئت علمی پژوهشگاه علوم انسانی و  
مطالعات فرهنگی

دکتر محمود بی‌جن‌خان عضو هیئت علمی دانشگاه تهران

دکتر اکبر حسابی عضو هیئت علمی دانشگاه اصفهان

دکتر علی رضاقلی فامیان عضو هیئت علمی دانشگاه پیام‌نور

دکتر مرتضی طاهری اردلی عضو هیئت علمی دانشگاه شهرکرد

دکتر مسعود قیومی عضو هیئت علمی پژوهشگاه علوم انسانی و  
مطالعات فرهنگی

دکتر علی رضاقلی فامیان عضو هیئت علمی دانشگاه پیام‌نور

دکتر مرتضی طاهری اردلی عضو هیئت علمی دانشگاه شهرکرد

دکتر سعیده ممتازی عضو هیئت علمی دانشگاه امیرکبیر

سید ابوالقاسم میرروشن‌دل عضو هیئت علمی دانشگاه گیلان

## هیئت داوران

دکتر سیدمصطفی عاصی عضو هیئت علمی پژوهشگاه علوم انسانی و

مطالعات فرهنگی

دکتر کامبیز بدیع عضو هیئت علمی موسسه تحقیقات

ارتباطات و فناوری اطلاعات

دکتر محمد بحرانی عضو هیئت علمی دانشگاه صنعتی شریف

دکتر محمود بی‌جن‌خان عضو هیئت علمی دانشگاه تهران

دکتر حسین صامتی عضو هیئت علمی دانشگاه صنعتی شریف

دکتر هشام فیلی عضو هیئت علمی دانشگاه تهران

دکتر اکبر حسابی عضو هیئت علمی دانشگاه اصفهان

دکتر مسعود قیومی عضو هیئت علمی پژوهشگاه علوم انسانی و

مطالعات فرهنگی

دکتر سعیده ممتازی عضو هیئت علمی دانشگاه امیرکبیر

دکتر علی رضاقلی‌فامیان عضو هیئت علمی دانشگاه پیام‌نور

دکتر بهروز مینایی عضو هیئت علمی دانشگاه علم و صنعت

ایران

دکتر محمدمهدی همایون‌پور عضو هیئت علمی دانشگاه صنعتی امیرکبیر

## کمیته اجرایی

زهرا ابراهیم‌بانکی، سهیلا ایزدی، محمدحسن ترابی، بیتا قوچانی، مریم محمدی،

ریحانه منظوری، طاهره همتی



## فهرست مطالب

- 
- ۹ پیشگفتار
- ۱۱ مسئله چندواژگی در پردازش نحو رایانشی زبان فارسی  
مسعود قیومی
- ۴۱ پیکره سازه: درختبانک بزرگ زبان فارسی در دستور سازه‌ای  
طباطبایی سیفی - صراف رضایی
- ۶۳ تولید درختبانک سازه‌ای زبان فارسی به روش نیمه خودکار  
محمدحسین دهقان و همکاران
- ۸۳ تجزیه سطحی معنایی جملات فارسی به کمک درخت  
ساخت سازه‌ای  
صغری لازمی و همکاران
- ۱۰۵ ارائه یک مدل بی‌نظمی بیشینه برای اصلاح خطای دستوری  
تطابق فعل و فاعل در زبان فارسی  
سیده زینب مفتاح - هشام فیلی
- ۱۲۷ معرفی سامانه استانداردساز و خطایاب متون علمی پژوهشگاه  
علوم و فناوری اطلاعات ایران  
ملوک‌السادات حسینی بهشتی - افتخارسادات هاشمی
- ۱۴۵ نظر کاوی خودکار نقد فیلم‌ها با رویکرد مقاوم‌سازی ماشین بردار  
پشتیبان  
امیرمحمود میر - جلال‌الدین نصیر
- ۱۶۵ سامانه خودکار خلاصه‌سازی با استفاده از روش تعبیه متن  
محمود کهنسال و همکاران
- ۱۸۷ شناسایی موجودیت‌های نامدار در شبکه‌های اجتماعی با  
رویکرد جمع‌سپاری  
شن‌آی بهراد و همکاران





## پیشگفتار

چهارمین همایش ملی زبان‌شناسی رایانشی در ۲۶ بهمن ۱۳۹۶ در پژوهشگاه علوم انسانی و مطالعات فرهنگی به همت انجمن زبان‌شناسی ایران و با همکاری نشر نویسه پاریس و پایگاه استنادی علوم جهان اسلام برگزار شد. تعداد ۲۲ مقاله از مراکز آموزشی و پژوهشی مختلف کشور به دبیرخانه همایش ارسال شد که مقالات ارسالی در چارچوب محورهای ذکرشده همایش بود. پس از داوری و برگزاری شورای کمیته داوران، با میزان پذیرش ۴۱ درصد، تعداد ۹ مقاله برای چاپ در مجموعه مقالات برگزیده شد که از این تعداد، ۸ مقاله به صورت سخنرانی ارائه شد.

زبان‌شناسی رایانشی میان‌رشته‌ای است که به بررسی مسائل زبانی با رویکرد رایانشی می‌پردازد. هدف اصلی این حوزه، درک و تولید زبان طبیعی در قالب گفتار و نوشتار است. برای رسیدن به این هدف، دانش زبانی باید به گونه‌ای الگوریتمی که قابلیت فهم و استفاده توسط رایانه را دارد، ارائه گردد. زبان‌شناسی رایانشی بخشی از علوم شناختی نیز محسوب می‌شود که با حوزه هوش مصنوعی به عنوان زیرمجموعه‌ای از علم رایانه، هم‌پوشی بسیار زیادی دارد. زبان‌شناسی رایانشی دو جنبه نظری و کاربردی دارد. جنبه نظری آن با زبان‌شناسی نظری و علوم شناختی در ارتباط است تا دانش زبانی برای درک و تولید به صورت صوری ارائه گردد. این گونه صوری‌سازی به شکلی دانش زبانی را شبیه‌سازی می‌کند که قابلیت پیاده‌سازی در نرم‌افزارهای رایانشی را داشته باشد. نرم‌افزارهای تهیه‌شده بیانگر جنبه کاربردی زبان‌شناسی رایانشی است.

مقالات پذیرفته‌شده در این همایش هر دو جنبه زبان‌شناسی رایانشی را پوشش داده‌است. تقریباً در همه مقالات، **زبان فارسی** به عنوان زبان هدف برای بررسی و پردازش انتخاب شده‌است. ویژگی رویکرد کاربردی‌سازی در حوزه زبان‌شناسی رایانشی و پردازش زبان فارسی سبب ارتقای جایگاه زبان فارسی در دنیای دیجیتال امروز می‌شود. واکاوی مقالات انتخاب‌شده مبین این نکته است که مقالات عمدتاً به معرفی یک الگوریتم یا روش با کاربرد مشخص برای بررسی

یک مسئله پرداخته‌است. خروجی این الگوریتم‌ها می‌تواند مستقیماً به تهیه نرم‌افزاری خاص برای پردازش زبان فارسی بیانجامد یا این که به‌عنوان بخشی از یک نرم‌افزار به کار رود. وجود این نرم‌افزارها می‌تواند به مطرح‌شدن سؤالات جدید و تلاش برای یافتن پاسخ به آنها در پژوهش‌های آینده منجر شود.

مقالات پذیرفته‌شده این همایش موضوعات زیر را در بر می‌گیرد: استانداردهای داده زبانی، تهیه درخت‌بانک سازه‌ای و تجزیه نحوی خودکار، خلاصه‌سازی خودکار متن، تشخیص موجودیت‌های نامدار، نظرکاوی، و تجزیه معنایی خودکار. بر خود لازم می‌دانم از تمامی بزرگانی که در شکل‌گیری این همایش تلاش کردند، تشکر و قدردانی نمایم. نخست از رئیس انجمن زبان‌شناسی ایران، سرکار خانم دکتر بلقیس روشن؛ از مسئولان محترم پژوهشگاه علوم انسانی و مطالعات فرهنگی؛ رئیس پژوهشکده زبان‌شناسی پژوهشگاه، جناب آقای دکتر مصطفی عاصی؛ دبیر اجرایی، سرکار خانم آذین شهریاری‌فرد و همچنین اعضای شورای اجرایی همایش که بزرگوارانه در اجرای همایش بسیار تلاش کردند؛ از سرکار خانم ریحانه منظوری که در ویرایش چکیده‌های انگلیسی کمک شایانی نمودند؛ و از جناب آقای امیر احمدی، مدیر محترم نشر نویسه‌پارسی، برای همراهی‌هایشان در اطلاع‌رسانی‌های همایش و همچنین چاپ مجموعه‌مقالات. در پایان شایسته است از تمامی اعضای محترم کمیته علمی و هیأت داوران همایش صمیمانه تشکر و قدردانی نمایم. همچنین از استادان، دانشجویان و پژوهشگران گرانقدری که با ارسال مقاله‌های پُر بار خود به غنای علمی این همایش افزودند و دستاوردهای پژوهشی خود را با علاقمندان به این حوزه در میان گذاشتند، صمیمانه سپاسگزاری می‌کنم.

مسعود قیومی

بهمن ۱۳۹۶

## مسئله چندواژگی در پردازش نحو رایانشی زبان فارسی

مسعود قیومی<sup>۱</sup>

### چکیده

این مقاله به بررسی چالش چندواژگی در پردازش نحو رایانشی زبان فارسی و ارائه راهکار برای رفع آن می‌پردازد. این چالش به دو دسته عمده تقسیم می‌شود: واحدهای واژگانی چند قطعه‌ای و چندقطعه‌ای‌های واژگانی یک واحدی. چالش دسته اول زمانی ظاهر می‌گردد که در یک زنجیره، چند واژه به‌اشتباه به یکدیگر جوش خورده‌اند، و یا واژه‌بست به میزبان ملحق شده‌است. در دسته دوم، چند زنجیره باید با هم ترکیب شوند تا یک واژه حاصل گردد. این دو چالش در پردازش نحوی و بن‌واژه‌سازی واژه‌ها متبلور است. برای رفع این دو دسته چالش، سه الگوریتم معرفی می‌شود که به‌ترتیب بر روی پیکره بی‌جن‌خان اجرا می‌گردند. ویژگی الگوریتم‌ها این است که در آنها از روش‌های قاعده‌مند و روش‌های مبتنی بر آمار استفاده شده‌است تا به‌طور منسجم بتوانند بر مشکلات حاصل از چندواژگی در پردازش نحو رایانشی زبان فارسی فائق آیند. پس از اعمال الگوریتم‌ها، کار ارزیابی با استفاده از داده‌آزمون نشان می‌دهد که با اجرای الگوریتم اول که در آن واژه‌بست از میزبان جدا می‌گردد، دقت ۸۰/۵۲ به‌دست می‌آید. با اجرای الگوریتم دوم برای رفع مشکل جوش خوردگی عناصر، دقت ۷۵/۴۳ به‌دست می‌آید. دقت ۸۶/۳۸ درصد با اجرای الگوریتم سوم برای ترکیب چند زنجیره و ساخت یک واژه به‌دست می‌آید.

**کلیدواژه‌ها:** چندواژگی، پردازش زبان طبیعی، نحو، پیکره زبانی.

<sup>۱</sup> استادیار گروه زبان‌شناسی، پژوهشگاه علوم انسانی و مطالعات فرهنگی.

## ۱. مقدمه

از جمله کاربردهای زبان‌شناسی پیکره‌ای یافتن مصادیق از داخل پیکره و استخراج آمار است. بنابراین داده زبانی باید به‌گونه‌ای معیار شده باشد تا امکان جستجوی دقیق و استخراج آمار دقیق فراهم گردد. کاربرد دیگر پیکره زبانی آن است که با نشانه‌گذاری داده زبانی امکان جستجوهای پیشرفته میسر گردد. این نشانه‌گذاری در سطوح مختلف زبانی از جمله نحو انجام می‌پذیرد. در انجام این امور، مؤلفه اصلی خط است، بنابراین هرگونه اشکالی در آن، بر داده، نتیجه جستجو، و آمار استخراج‌شده تأثیرگذار است. خط فارسی به‌دلیل ویژگی‌هایی که دارد سبب می‌شود پردازش رایانشی داده زبانی را با اشکال مواجه کند. آنچه در این مقاله به آن پرداخته می‌شود چندواژگی به‌عنوان یک چالش اصلی در پردازش رایانه‌ای زبان فارسی به‌خصوص در قسمت نحو است. علی‌رغم پژوهش‌های انجام‌شده در این زمینه مانند شریفی‌آتشگاه و بی‌جن‌خان (۲۰۰۹) و طباطبایی و صراف (۱۳۹۶)، نیاز است با نگاهی عمیق‌تر و اشراف بیشتر این مسئله مورد بررسی و ارزیابی قرار گیرد. این مقاله در ۷ بخش تهیه شده‌است. در بخش ۲، چالش‌هایی که برای تهیه یک پیکره برای زبان فارسی با آن مواجه می‌شویم، بیان می‌شود. بخش ۳ به معرفی داده و ابزار به‌کاررفته در این پژوهش می‌پردازد. در بخش ۴، چندواژگی به‌عنوان یکی از چالش‌ها مطرح شده و توضیح داده می‌شود. بخش ۵ به معرفی سه الگوریتم و ارائه راهکارهایی برای حل چالش چندواژگی می‌پردازد و توضیح داده می‌شود. در بخش ۶، نتایج به‌دست‌آمده بحث و بررسی خواهد شد. در انتها در بخش ۷، به نتیجه‌گیری راهکارهای ارائه‌شده خواهیم پرداخت.

## ۲. چالش‌های تهیه پیکره زبان فارسی

قیومی و ممتازی (۲۰۰۹)، قیومی و همکاران (۲۰۱۰)، و شمس‌فرد (۲۰۱۱) چالش‌هایی را که به‌هنگام تهیه پیکره زبان فارسی با آنها مواجه می‌شویم، بررسی کرده‌اند که در این بخش به‌طور خلاصه به آنها اشاره می‌شود. کدگذاری حروف فارسی و عربی و تداخل بین کد حروف «ی» و «ک» فارسی با کد عربی حروف «ی» و «ک» یکی از این چالش‌هاست. این تداخل در تهیه فهرست الفبایی نیز تأثیرگذار است. چالش دیگر چندواژگی است که در ادامه این بخش توضیح داده خواهد شد. سبک نوشتاری از نظر نوشتار معیار، فرامعیار، یا فرومعیار چالش دیگر است که باید در تهیه پیکره معیار که به‌عنوان نماینده زبان طلقی می‌شود در نظر گرفته شود. خلایقیت زبانی متجلی‌شده در خط نیز چالش دیگر است. این ویژگی در پیامک‌ها و نظرات کاربران و وبلاگ‌ها کاملاً مشهود است. ابهام به‌دلیل هم‌نویسی واژه‌ها چالشی دیگر است. ممکن است چندواژگی سبب هم‌نویسی گردد و سپس چالش ایجاد گردد. عامل دیگر هم‌نویسی مانند عدم‌درج یا تغییر حروف و علائم قرضی از خط عربی مانند همزه یا تنوین سبب ایجاد هم‌نویسی شده و پردازش خط را با چالش مواجه کند. متأسفانه علی‌رغم وجود دستور املائی مصوب فرهنگستان زبان و ادب فارسی (۱۳۸۹) فهرستی از استثناها مطرح می‌گردد که این خود منجر به آشفتگی و عدم قطعیت در املاء و پراگندگی آمار واژه می‌گردد. عدم وجود هرگونه معیار برای نگارش واژه‌های لاتین با حروف فارسی، مانند اسم لاتین نویسنده، در یک متن فارسی منجر به خلق واژه‌های جدید با توزیع آماری مستقل می‌گردد، بنابراین این خود یک چالش محسوب می‌شود.

تعدادی از این چالش‌ها را قیومی و همکاران (۱۳۹۴) تحت عنوان تنوع نگارشی بررسی کرده و با معرفی یک الگوریتم تلاش کرده‌اند، تنوع‌های نگارشی را به‌طور خودکار بیابند. در این مقاله چندواژگی و رفع مشکل آن

مطرح می‌گردد تا داده حاصل را بتوان برای نشانه‌گذاری برای سطوح دیگر زبان‌شناسی استفاده کرد.

شریفی‌آتشگاه و بی‌جن‌خان (۲۰۰۹) به‌صورت پیکره‌بنیان ویژگی چندواژگی را در زبان فارسی مطرح کرده‌اند. در بررسی آنها چندواژگی به دو دسته تقسیم می‌گردد. این دو دسته، محصول عدم‌رعایت فاصله‌گذاری صحیح بین واحدهای چندواژه‌ای است. دسته اول، واحدهای واژگانی چند قطعه‌ای هستند که چندواژه با هم یک واحد واژگانی می‌سازند. در این دسته، فاصله کامل<sup>۲</sup> به جای فاصله مجازی<sup>۳</sup> درج شده‌است. بنابراین بین عناصر واژگانی فاصله کامل قرار گرفته‌است و با توجه به این که فاصله کامل به‌عنوان مرز واژه تلقی می‌گردد، سبب می‌شود یک واحد واژگانی به‌صورت چند واژه توسط رایانه تشخیص داده شود. برای مثال، عدم فاصله‌گذاری صحیح منجر می‌شود واژه «درغیراین‌صورت» به‌صورت «درغیراین صورت» نگارش گردد و دو واژه «درغیراین» و «صورت» توسط رایانه تشخیص داده شود. واژه‌های مرکب، مانند «برای اینکه»، و یا مواردی که بین وند تصریفی یا اشتقاقی و شکل پایه‌ای واژه، فاصله کامل به جای فاصله مجازی درج شده است، مانند «دانش آموز»، «بی‌هدف»، «فروشگاه‌ها» و غیره را می‌توان در این دسته قرار داد.

دسته دوم، چندقطعه‌ای‌های واژگانی یک واحدی هستند که چند واحد واژگانی با هم یک واژه را تشکیل می‌دهد. در این دسته، عدم‌درج فاصله کامل سبب می‌شود چندین واژه به‌عنوان یک واژه تلقی گردد. منشاء اصلی این مشکل این است که تعدادی از حروف فارسی، مانند «آ»، «ا»، «د»، «ذ»، «ر»، «ز»، «ژ»، «و»، شکل متصل به حرف بعدی ندارند؛ بنابراین اتصال این حروف به حرف بعدی اختلالی را در امر خواندن ایجاد نمی‌کند. برای مثال عدم‌درج فاصله سبب می‌شود واحد «و یا در برابر او ایستادم» به‌صورت یک

<sup>2</sup>white space

<sup>3</sup>pseudo-space

واحد «ویادربرابراواستادم» نگارش شود که این زنجیره به صورت یک واحد توسط رایانه تشخیص داده می‌شود. عدم رعایت فاصله‌گذاری ممکن است به ابهام بیانجامد و مشخص نگردد آیا این زنجیره از حروف باید از منظر یک واحد چندواژه‌ای تحلیل گردد یا از منظر یک واژه چندواحدی یا خود یک واژه بسیط و مستقل است. برای مثال «وبا» می‌تواند یک نوع بیماری باشد؛ بنابراین یک واژه بسیط است. همچنین ممکن است یک واحد چند واژه‌ای «و» و «با» باشد. مثال‌های دیگری چون «مادر»، «توهم»، «راداری»، «بریزید»، و غیره بیانگر این ابهام است که نشان می‌دهد تصحیح داده به صورت خودکار نمی‌تواند کار ساده‌ای باشد. حالت کوتاه‌شدگی واژه‌ها در متون ادبی، مانند «مرا»، «زوی»، «کزو» و غیره، و یا در گفتار، مانند «بچته»، «چته» یا «بچتو» را می‌توان در دسته یک واحد چندواژه‌ای قرار داد. ضمائر ملکی یا عنصر فعلی که به صورت واژه‌بست به پایه متصل می‌شوند، مانند «کتابش»، «درعذابند»، یا «داراست»، در این دسته قرار می‌گیرند.

به‌هنگام کار با پیکره و نشانه‌گذاری داده در سطوح مختلف زبان‌شناسی دو چالش مطرح شده می‌تواند مشکل‌زا باشد. در پیکره بی‌جن‌خان (بی‌جن‌خان، ۱۳۸۳) مقوله دستوری واژه‌ها مشخص شده‌است که در بخش بعدی توضیح داده خواهد شد. مجموعه برچسب‌های به‌کارگرفته‌شده برای نشانه‌گذاری این داده به‌گونه‌ای تعریف شده‌است که می‌تواند مشکل دسته اول را توجیه نماید ولی در مورد مشکل دسته دوم، مقوله دستوری واژه‌ها به صورت مستقل مشخص شده‌است و مقوله دستوری عناصر ترکیب‌شده با هم نامشخص است. این مشکلات در سطح بالاتر تحلیل زبان‌شناختی مانند نحو که در آن باید ترسیم نمودار درختی جمله از روابط واژه‌ها با هم و نه مستقل از هم به‌دست آید، به‌خوبی خود را نشان می‌دهد. این مشکلات در بن‌واژه‌سازی نیز اشکالاتی ایجاد می‌کند. در این مقاله تلاش می‌شود برای رفع هر یک از این دو دسته چالش‌ها، راهکاری معرفی گردد تا قابلیت استفاده در پردازش نحو

جملات فارسی و همچنین تهیۀ دادگان درختی فارسی و بن‌واژه‌سازی داشته باشد.

### ۳. معرفی داده و ابزار به‌کاررفته

در پیکرۀ بی‌جن‌خان (۱۳۸۳) مقولات دستوری به‌صورت سلسله‌مراتبی برچسب‌گذاری شده‌است (بی‌جن‌خان و همکاران، ۲۰۱۱). در این برچسب‌گذاری از استاندارد ایگلز<sup>۴</sup> (لیچ و ویلسون، ۱۹۹۹) استفاده شده‌است. نحوه سلسله‌مراتبی این برچسب‌گذاری به‌این‌صورت است که ابتدا مقولۀ دستوری اصلی واژه مشخص می‌گردد و سپس ویژگی‌های صرفی-نحوی و تا حدی معنایی برای هر واژه مشخص می‌شود. در جدول ۱ این نحوه برچسب‌گذاری نمایش داده شده‌است. ترتیب قرارگرفتن اطلاعات صرفی-نحوی منجر به ایجاد ۵۸۶ برچسب متنوع شده‌است.

در این برچسب‌گذاری از ویژگی‌های معنایی برای ابهام‌زدایی هم‌نویسه‌ها استفاده شده‌است. بنابراین در نمونه‌های جدول ۱، اضافه‌شدن ویژگی «مکان» سبب می‌شود «دفتر» به‌معنای «محل کار» از «دفتر» به‌معنای «دفتر یادداشت» مجزا گردد.

جدول ۱: نمونه برچسب‌گذاری در پیکرۀ بی‌جن‌خان

| واژه | برچسب           | تفسیر برچسب                 |
|------|-----------------|-----------------------------|
| دفتر | N,COM,SG,LOC    | اسم، عام، مفرد، مکان        |
| دفتر | N,COM,SG,LOC,EZ | اسم، عام، مفرد، مکان، اضافه |
| دفتر | N,COM,SG        | اسم، عام، مفرد              |
| دفتر | N,COM,SG,EZ     | اسم، عام، مفرد، اضافه       |

<sup>4</sup> EAGLES



در بخش دوم مقاله حاضر دو دسته چندواژگی مطرح شد. در دسته اول چندواژگی، چند واحد واژگانی به عنوان یک واحد تلقی می شود. کوتاه شدگی و عدم رعایت فاصله گذاری سبب می شود بین عناصر واژگانی جوش خوردگی اتفاق افتد و این عناصر به شکل یک واحد تلقی گردد. واژه بست نیز عنصر دیگری است که ذاتاً نیاز به عنصری دیگر برای جوش خوردگی دارد. در اینجا به نمونه هایی از پیکره بی جن خان و نحوه برچسب دهی آنها اشاره خواهد شد که در جدول شماره ۲ نمایش داده شده است. در بخش ۴، این نحوه برچسب دهی مقولات دستوری به عنوان چالش چندواژگی در تهیه دادگان درختی مطرح می شود.

جدول ۲: برچسب گذاری چندواژگی به صورت یک واحد در پیکره بی جن خان

| واژه   | برچسب                  | تفسیر برچسب   |
|--------|------------------------|---|
| وبا    | CONJ,P<br>N,COM,SING   | حرف ربط، حرف نشانه را<br>اسم، عام، مفرد                         |
| بچتو   | N,COM,SING,2,POSTP     | اسم، عام، مفرد، واژه بست<br>ضمیری دوم شخص مفرد،<br>حرف نشانه را |
| کتابتو | N,COM,SING,2, POSTP    | اسم، عام، مفرد، واژه بست<br>ضمیریدوم شخص مفرد،<br>حرف نشانه را  |
| کزو    | CONJ,P,PRO,PERS,SING,3 | حرف ربط، حرف اضافه، ضمیر<br>شخصی سوم شخص مفرد                   |
| اویم   | V,COP,PROC,1           | فعل، ضمیر پیش فعلی، اول<br>شخص مفرد                             |

| واژه     | برچسب          | تفسیر برچسب   |
|----------|----------------|---|
| درعذابند | V,COP,PREPNC,6 | فعل، حرف‌افزافه به‌همراه گروه اسمی متممی، سوم شخص جمع |
| جدیدش    | AJ,SIM,3       | صفت، ساده، واژه‌بست ضمیری سوم شخص مفرد                |

در دسته دوم، یک واحد واژگانی از ترکیب یک یا چند واژه حاصل می‌گردد. برای این ترکیب نیاز است واژه‌های موردنظر با هم ترکیب شوند تا یک واحد واژگانی ساخته شود. برای انجام این کار می‌توان از ابزاری مانند سامانه کلاک<sup>۵</sup> (سیمو و اوسنوا، ۲۰۰۲) استفاده کرد که شاخص‌ترین ویژگی‌های این ابزار در ادامه توضیح داده خواهد شد.

سامانه کلاک<sup>۶</sup> که به‌طور رایگان قابل دسترس است، براساس XML بنا نهاده شده‌است. این سامانه مجهز به زبان XPATH است تا جستجو در اسناد را میسر سازد. در این سامانه، قواعد دستوری می‌تواند با استفاده از اطلاعات واژی-نحوی به‌صورت عبارات باقاعده<sup>۷</sup> بیان گردد. برای جلوگیری از فرازایش<sup>۸</sup> این قواعد، بافت محلی قواعد مورد توجه قرار می‌گیرد تا با تعریف محدودیت‌هایی برای قواعد، از فرازایش این قواعد کاسته شود. قواعد تعریف‌شده در این سامانه به‌عنوان «ماشین‌های حالت محدود»<sup>۹</sup> نقش آفرینی می‌کند که براساس دستور آبخاری<sup>۱۰</sup> آبنی (آبنی، ۱۹۹۶) این قواعد به‌صورت

<sup>5</sup>CLARK

<sup>6</sup> <http://www.bultreebank.org/clark/>

<sup>7</sup> Regular Expression

<sup>8</sup> Over-generate

<sup>9</sup> Finite State Automata

<sup>10</sup> Cascaded Grammar

سلسله‌مراتبی با هم در ارتباط بوده و خروجی حاصل از یک قاعده (ماشین) به‌عنوان ورودی قاعدهٔ بعدی (ماشین بعدی) تعریف می‌گردد. همان‌طور که پیشتر عنوان شده بود، ساختار داده در کلارک براساس XML است؛ بنابراین هر واژه در یک گره قرار می‌گیرد و اطلاعات زبان‌شناختی مربوط به آن واژه به‌صورت مشخصه<sup>۱۱</sup> و ارزش<sup>۱۲</sup> بیان می‌شود. دسترسی به اطلاعاتی مانند صورت واژه و مقولهٔ دستوری آن، به انعطاف‌پذیری در تعریف قواعد باقاعده و همچنین محدودکردن اعمال این قواعد می‌افزاید. در کلارک سه علامت مفهوم خاصی دارند که برای افزایش انعطاف‌پذیری قواعد می‌توان از این علائم به‌هنگام نگارش قواعد استفاده کرد. علامت @ به‌مفهوم هیچ یا یک حرف است، علامت % به‌مفهوم هیچ یا بیش از یک حرف است، و علامت # به‌مفهوم یک یا بیش از یک حرف است.

#### ۴. چالش چندواژگی در تهیهٔ دادگان درختی

در بخش ۲، ویژگی‌های برجسته‌های دستوری واژه‌ها معرفی شد. آنچه در برجسته‌ها مشاهده می‌شود، این است که تلاش شده‌است برجسته‌ها با شکل نوشتاری منطبق گردد که در این صورت نحوهٔ تعریف برجسته‌ها در راستای خط فارسی بوده و شیوهٔ نگارش در برجسته‌ها لحاظ شده‌است. اتخاذ چنین روشی در تحلیل و نمایش اطلاعات زبان‌شناختی می‌تواند در سطح مشخص شدن نقش دستوری واژه مفید باشد، ولی برای سطوح تحلیل عمیق‌تر مانند تهیهٔ تجزیهٔ درختی جمله ناکارآمد است. از آنجاکه دادگان درختی<sup>۱۳</sup> جملات می‌تواند برای آموزش یک تجزیه‌گر خودکار<sup>۱۴</sup> توسط رایانه مورد استفاده قرار گیرد، واضح و مشخص بودن اطلاعات زبان‌شناختی باید در ساختار داده

<sup>11</sup>attribute

<sup>12</sup>value

<sup>13</sup>treebank

<sup>14</sup>parsing

لحاظ گردد تا بهترین سطح کارایی تجزیه‌گر به‌دست آید. برای درک بهتر این چالش، جمله (۱) را در نظر بگیرید.

(۱) حسین آمد و با حسن به مدرسه رفت.

«با» هسته گروه حرف‌اضافه‌ای عبارت «با حسن» است. در قواعد ساخت سازه‌ای تعریف‌شده در کلاک، حرف‌اضافه و اسم با هم ترکیب می‌شوند تا یک گروه حرف‌اضافه‌ای بسازد. جوش‌خوردگی «و» و «با» و درج برجسب دستوری «P, CONJ» در پیکره سبب می‌شود امکان اتصال دو واژه «با» و «حسن» میسر نگردد. دلیل آن این است که در زنجیره «و با حسن» هسته حرف‌اضافه‌ای و برجسب آن مشخص نیست تا قاعده تعریف‌شده اعمال گردد. در مثال (۲) تا (۴) پدیده واژه‌بست مطرح است.

(۲) کتابتو از حسن گرفت.

(۳) کتاب جدیدش را از حسن گرفت.

(۴) آنها سخت درعذابند.

در «کتابتو» دو واژه‌بست وجود دارد: یکی نقش ضمیر ملکی را داراست، و دیگری نقش‌نمای مفعولی. می‌دانیم ضمیر ملکی و حرف نشانه «را» در چارچوب دستور ساخت سازه‌ای هسته‌بنیان<sup>۱۵</sup> متمم هسته اسم (سمولیان و تسنگ، ۲۰۱۰؛ مولر و قیومی، ۲۰۱۰) تلقی می‌گردند که در مثال (۲) این دو واژه‌بست متمم «کتاب» هستند. از آنجا که این عناصر واژگانی (واژه‌بست‌ها) با نقش‌های دستوری متفاوت (واژه‌بست ملکی و واژه‌بست نقش‌نما) به‌صورت یک واحد آمده‌است، نیاز است از پایه منفک گردند تا نقش و روابط دستوری آنها به‌عنوان یک واحد واژگانی کاملاً مستقل مشخص گردد و تحلیل صحیح جمله به‌دست آید.

<sup>15</sup> Head-driven Phrase Structure Grammar