

# روشی نوین جهت صرف واژه های فارسی

بهرنگ قاسمی زاده سعید رحیمی  
حیدر سمیاری عباس کوچاری  
مرتضی سالاریان مجید نم نبات  
علی ترکمنی لقمان براری

دپارتمان بومی سازی، شرکت دیجیتال کلون  
[QasemiZadeh@DigitalClone.net](mailto:QasemiZadeh@DigitalClone.net)

به طور مستقل می تواند در جمله ظاهر شود. (۲) تکواژ مقید<sup>۵</sup> مانند "ی" در "شادی" که به تنهایی نمی تواند در جمله به کار رود. تکواژ مقید خود بر دو گونه است: (۱) تکواژ صرفی<sup>۶</sup> مانند شناسه های فعلی (ام، ای، ...) که افزودن آنها به تکواژ پایه به دلیل الزامات صرفی جهت ایفای نقش نحوی است و سبب تغییر مقوله زبانی واژه و نیز معنا نمی شود. در زبان فارسی این تکواژها بیشتر به صورت پسوند ظاهر می شوند. (۲) تکواژ اشتقاقی<sup>۷</sup> مانند "گاه" در "دانشگاه" که معمولاً حضور آن در واژه سبب تغییر مقوله واژه شده و تغییر معنایی را به همراه دارد. این تکواژها هم به صورت پیشوند و هم به صورت پسوند ظاهر می شوند.

از آنجا که تکواژهای مقید به پیش و یا پس یک تکواژ مقید و یا آزاد دیگر افزوده می شوند، به آنها وند<sup>۸</sup> می گویند. به فرایند اضافه کردن وند به تکواژ پایه وندافزایی<sup>۹</sup> گویند. اگر این وند در پیش واژه قرار گیرد به آن پیشوند<sup>۱۰</sup>، و اگر در پس واژه قرار گیرد به آن پسوند<sup>۱۱</sup> گویند. بنابراین وند تکواژ مقیدی است که می توان آن را به دو نوع تصریفی و اشتقاقی تقسیم کرد.

اگر مقایسه ای ساده بین وند تصریفی و اشتقاقی صورت گیرد، خواهیم دید که وند تصریفی نقش نحوی دارد و واژه را برای ایفای نقش معینی در ساخت نحوی پردازش می کند و هرگز مقوله زبانی<sup>۱۱</sup> تکواژ پایه را تغییر نمی دهد. کاربرد وند تصریفی قیاسی است به این مفهوم که با همه اعضای یک مقوله از واژه ها می تواند به کار رود و کمتر استثنا می پذیرد. اما وندهای اشتقاقی نقش واژه سازی داشته و واژه جدید می سازد. وند اشتقاقی اغلب مقوله زبانی تکواژ پایه را تغییر می دهد. کاربرد وندهای اشتقاقی سماعی است یعنی با همه اعضای

**چکیده:** در این مقاله رویکردی نوین برای صرف واژه در زبان فارسی ارائه شده است. صرف بخشی از دستور زبان است که ساخت واژه را مورد بررسی قرار می دهد از مهمترین مباحث صرف، تصریف و اشتقاق است که در زبان فارسی به کمک وندافزایی صورت می گیرد. در این مقاله روشی جدید با پیچیدگی محاسباتی مناسب جهت حل مشکل صرف تصریفی واژه ها در زبان فارسی ارائه شده است. روش پیشنهادی به سادگی قابل بسط برای دیگر زبان های هند و اروپایی است که در آن تصریف واژه از طریق وندافزایی صورت می گیرد. در روش پیشنهادی، مسئله تحلیل لغوی به مسئله جستجوی جزئی هدایت شده بر روی درخت جستجوی سه تایی تبدیل شده است. مشکل موجود در صرف تصریفی شامل استخراج قوانین، تعداد زیاد قوانین و نحوه اعمال آن بر روی واژه های زبان است. از جمله مشکلاتی که در تحلیل صرفی، مبتنی بر دانش و قوانین با آن مواجه هستیم می توان از تعداد زیاد قوانینی که می توانند همزمان اعمال شوند، عدم یکنواختی، و فضای جستجوی بزرگ نام برد. در این پژوهش سعی شده است روشی برای حل این مشکلات ارائه شود.

**کلمات کلیدی:** صرف واژه های فارسی، پردازش زبان طبیعی، تحلیل لغوی، واژگان، بازنمایی دانش

## ۱- مقدمه

می دانیم واحد زبان، جمله است و جمله خود از واحدهای کوچکتری که واژه<sup>۱</sup> نام دارد تشکیل شده است. از کنار هم نشستن واژه ها با توجه به نحو زبان، جمله وار و جمله ساخته می شود. واژه از اجزای کوچکتری به نام تکواژ<sup>۲</sup> تشکیل شده است. به بررسی ساختمان واژه، صرف<sup>۳</sup> (ساخت واژه) گویند. تکواژ که کوچکترین واحد معنادار زبان است به دو دسته تقسیم می شود: (۱) تکواژ آزاد<sup>۴</sup> مانند "درخت" که

<sup>5</sup> Bound Morpheme

<sup>6</sup> Inflectional Morpheme

<sup>7</sup> Derivational Morpheme

<sup>8</sup> Affix

<sup>9</sup> Prefix

<sup>10</sup> Suffix

<sup>11</sup> Part of Speech (POS)

<sup>1</sup> Word

<sup>2</sup> Morpheme

<sup>3</sup> Morphology

<sup>4</sup> Free Morpheme

آن مشخص شده است آموزش داده شود. متاسفانه چنین پیکره‌هایی برای استفاده عمومی در زبان فارسی منتشر نشده است. در [۹] یک سیستم تحلیلگر لغوی غیرقطعی<sup>۵</sup> مبتنی بر آموزش حافظه مدار<sup>۶</sup> ارائه شده است. در این روش یک نگاهت مستقیم از حروف در یک متن به لایه‌های متفاوتی که اطلاعات صرفی را در خود رمز نموده‌اند، تعریف شده است. این سیستم نیز پیش از استفاده نیاز است آموزش داده شود. در حقیقت این سیستم با تشخیص مرزبندی‌ها در یک واژه، آنها را طبقه بندی و برچسب دهی می‌نماید.

در [۱۲] یک روش مبتنی بر واژگان برای حل مشکل تحلیل لغوی ارائه شده است. در این مقاله از یک گراف غیرمدور برای ذخیره اطلاعات واژگان به همراه اطلاعات ساخت واژه استفاده شده است که پردازش فقط از طریق یکبار جستجو در این گراف صورت می‌گیرد.

برای زبان فارسی تا به حال چندین کار انجام شده است به‌عنوان مثال [۷] و [۱۳]. در [۱۳] از یک معماری مبتنی بر ماشین حالت متناهی جهت ارائه یک تحلیل گر لغوی استفاده شده است. آقای ریاضی در [۷]، یک سیستم تحلیلگر لغوی دو سطحی برای تحلیل صرفی و اشتقاقی زبان فارسی ارائه نموده است. مدل Perslex با الهام از مدل Englex [۱۴] ارائه شده است. مشابه چنین کاری در [۱۵] و به کمک ابزارهای Xerox finite state برای ارائه یک سیستم تحلیلگر لغوی برای زبان فارسی استفاده شده است.

هدف تحلیلگر لغوی، تحلیل ویژگی‌های صرفی واژه و تشخیص تغییرات صورت گرفته در آن به واسطه نقش نحوی آن است. برای رسیدن به این هدف، ابتدا یک دسته‌بندی جدید از مقوله‌های زبانی<sup>۷</sup> در زبان فارسی ارائه شده است. سپس با توجه به این دسته‌بندی و استفاده از روشی مناسب جهت بازنمایی واژگان، یک روش جدید برای تحلیل صرفی واژه‌ها ارائه شده است.

در بخش ۲، مقوله‌های زبانی در فارسی شرح داده می‌شود. در ادامه و در بخش ۳ به واژگان، ساختار و رابطه آن با تحلیل لغوی اشاره می‌شود. در بخش ۴ روش پیشنهادی در طراحی تحلیلگر لغوی آورده شده است. ساختار سیستم در بخش ۵ شرح داده می‌شود. بخش ۶ بیانگر نتایج به‌دست آمده است و در انتها به جمع بندی و پیشنهاد برای کارهای آینده پرداخته می‌شود.

## ۲- مقوله‌های واژگانی در زبان فارسی

به‌طور کلی دسته‌بندی‌های انجام شده بر روی مقوله‌های زبانی در زبان فارسی به دو دسته تقسیم می‌شود. روش سنتی و روش نوین. روش سنتی بر پایه دیدگاه‌های قدیمی ساخت زبان استوار است. در

یک مقوله از واژه‌ها به کار نمی‌رود. تعداد وندهای تصریفی به مراتب کمتر از وندهای اشتقاقی است. بر همین اساس، بسامد کاربرد ونده تصریفی از ونده اشتقاقی بسیار بالاتر است. از نقطه نظر محل قرار گرفتن وندهای اشتقاقی در مقایسه با ونده تصریفی، معمولاً ونده اشتقاقی به پایه نزدیکتر است و در درونترین لایه ساخت‌واژه قرار می‌گیرد. در حالی که ونده تصریفی در برونترین لایه ساخت‌واژه قرار می‌گیرد. اما باید توجه نمود که در زبان فارسی بعضی از وندهای تصریفی قبل از وندهای اشتقاقی قرار می‌گیرند. ونده تصریفی و اشتقاقی را می‌توان از لحاظ روابط معنایی از هم بازشناخت. در ونده تصریفی رابطه میان معنی تکواژه پایه و معنی تکواژ جدید حاصل از افزایش ونده به پایه کاملاً قابل پیش‌بینی است و از قوانین مشخصی تبعیت می‌نماید؛ در حالی که در ونده اشتقاقی رابطه میان تکواژ پایه و معنی تکواژ جدید بدست آمده از ونده‌افزایی غیرقابل پیش‌بینی است [۱].

اگر وندها را با توجه به محل قرارگیریشان نسبت به پایه دسته‌بندی کنیم، سه دسته ونده خواهیم داشت: پیشوند، پسوند و میانوند<sup>۱</sup>. لازم به تذکر است که در فارسی میانوند وجود ندارد [۱۱][۲]. یکی از مشکلات مهم در زبان فارسی که باید به آن توجه نمود، هم‌نویسه بودن گروهی از وندهای تصریفی با وندهای اشتقاقی است؛ چراکه در تحلیل واژه‌های حاصل از افزایش چنین وندهایی، با ابهام روبرو هستیم مانند "ی" نشانه نکره در مقایسه با "ی" ونده اشتقاقی. آنچه در این مقاله مورد بحث قرار می‌گیرد، بررسی ساخت‌واژه واژه‌های حاصل از افزایش وندهای مقید تصریفی است.

مسئله تحلیل لغوی یکی از مسائل قدیمی در پردازش زبان طبیعی می‌باشد. امروزه برای اکثر زبان‌های زنده دنیا، سیستم‌های پردازش لغوی وجود دارد. پوشش مناسب از واژه‌های زبان، بار محاسباتی کم، سرعت پردازش بالا و همچنین دقت کافی از مهمترین مسائل در طراحی پردازشگرهای لغوی است. روش‌های متعددی برای پردازش لغوی پیشنهاد شده است. بیشتر کارهای انجام شده، طراحی تحلیلگرهای لغوی مبتنی بر ماشین حالت متناهی<sup>۲</sup> بوده است [۳][۴][۵] [۶][۷]. علاوه بر آن روش‌های آماری و یا مبتنی بر یادگیری ماشین نیز برای استفاده در طراحی سیستم‌های ریشه‌یاب<sup>۳</sup> خودکار ارائه شده است [۸][۹][۱۰][۱۱].

در [۸] یک معماری مستقل از زبان مبتنی بر مدل مخفی مارکوف برای تحلیلگر لغوی ارائه شده است. این سیستم برای زبان‌های چینی، ژاپنی و انگلیسی آزمایش شده است. سیستم فوق پیش از استفاده نیاز است تا بر روی پیکره<sup>۴</sup> زبانی که مقوله زبانی واژه‌های

<sup>۱</sup> Infix

<sup>۲</sup> Finite State Automata

<sup>۳</sup> Stemmer

<sup>۴</sup> Corpus

<sup>۵</sup> Non Deterministic

<sup>۶</sup> Memory Based Learning

<sup>۷</sup> Part of Speech (POS)

این دسته‌بندی، اجزای کلام فارسی به هفت گروه تقسیم بندی شده است [۱۶]. از آنجا که روش سنتی بر پایه دیدگاه‌های قدیمی استوار است، چندان به رابطه صرف و نحو تکیه نداشته و بیشتر به تحلیل واژه مستقل از روابط نحوی می‌پردازد [۱۷][۱۸] دسته‌های ارائه شده در روش سنتی عبارتند از [۱۶]:

- اسم: می‌تواند مستقیماً و مستقلاً نهاد(مسند الیه) جمله باشد و آن برای نامیدن شخصی یا حیوانی یا چیزی و یا مفهومی به کار می‌رود.
- صفت: حالت و مقدار و شماره یا یکی دیگر از چگونگی‌های اسم را می‌رساند.
- فعل: در جمله جایگاه اسناد را اشغال می‌کند، یعنی یا خود به نهاد اسناد داده می‌شود یا کلمه ای را به نهاد اسناد می‌دهد و به تنهایی یا به کمک وابسته‌هایی اغلب به چهار مفهوم دلالت می‌کند: ۱. شخص ۲. شمار ۳. زمان ۴. حالت یا کار.
- قید: کلمه یا گروهی که مفهومی به مفهوم فعل و یا صفت یا مسند یا قید و مصدر دیگر می‌افزاید.
- ضمیر: کلمه ای که معمولاً به جای اسم می‌نشیند.
- عدد
- اصوات
- حروف: نقش دستوری واژه‌های دیگر را نشان می‌دهد.

این دسته‌بندی دو مشکل اساسی دارد. اول این که دسته‌بندی سنتی اطلاعات کافی را جهت تحلیل نحوی زبان به کمک روش‌های نوین زبان‌شناسی فراهم نمی‌آورد. مشکل دوم این دسته‌بندی، عدم وجود دسته و گروهی خاص برای برخی از اجزای کلام در زبان فارسی است. این همان چیزی است که در زبان‌شناسی سنتی با نام ادات شناخته می‌شود. اما وجود چنین دسته‌بندی سبب فراهم آوردن اطلاعات افزوده‌ای نخواهد بود و کمکی در جهت تسهیل تحلیل رایانه‌ای زبان فراهم نمی‌آورد.

روش‌های نوین زبان‌شناسی تعاریف جدیدی از مقوله‌های زبانی یا اجزای کلام در زبان فارسی ارائه داده‌اند [۱۹][۲۰]. با توجه به نواقص موجود در دسته‌بندی سنتی، [۲۰] دسته‌بندی جدیدی از مقوله‌های زبانی با توجه به زبان‌شناسی نوین ارائه کرده است. در این دسته بندی تکواژها بر پایه ویژگی‌های معنایی، صرفی و نحوی در چهار دسته جای می‌گیرند. در این دسته‌بندی تکواژ به‌عنوان کوچکترین واحد زبان بررسی می‌شود. دسته‌بندی مقوله‌های زبانی واژه های فارسی مطابق با این نظریه عبارتست از:

- تکواژ واژگانی یا قاموسی
  - ۱- پایه فعل ۲- اسم ۳- صفت ۴- قید
  - ۲- واژه‌های دستوری یا نقشی
  - ۳- حرف اضافه ۲- حرف ربط ۳- ضمیر

- تکواژ اشتقاقی
- تکواژهای تصریفی

- ۱- نشانه‌های جمع
- ۲- نشانه‌های برتر برای واژه‌های صفت
- ۳- نشانه برتر برای واژه‌های قید
- ۴- نشانه‌های ماضی ساز
- ۵- نشانه‌های صفت مفعولی
- ۶- جزء پیشین فعلی
- ۷- شناسه‌های فعلی
- ۸- پی‌بست‌های ماضی نقلی

- واژه‌بست

- ۱- کسره اضافه و وصفی
- ۲- صورت‌های بودن
- ۳- ی نکره
- ۴- ضمائر متصل
- ۵- الفِ ندا : خدایا

این دسته‌بندی یک گروه‌بندی دقیق از مقوله‌های زبانی در فارسی ارائه می‌دهد؛ اما در بعضی موارد فاقد ارزش و اعتبار کافی برای تحلیل رایانه‌ای صورت نوشتاری واژه‌ها در زبان فارسی است. نگارندگان بر این عقیده‌اند که چگونگی دسته‌بندی ارائه شده برای مقوله‌های زبانی مستقیماً می‌تواند بر روی الگوریتم‌های طراحی شده تأثیرگذار باشد. به عبارتی دیگر می‌توان قسمتی از دانش زبانی مورد نیاز جهت تحلیل زبان را با کمک دسته‌بندی‌ها و روابط منطقی که میان این دسته‌ها برقرار است، به‌صورت ضمنی ارائه کرد. با توجه به آنچه که گفته شد و با توجه به مطالعات صورت‌گرفته و تجربیات کسب‌شده در طراحی و پیاده‌سازی نسخه‌های اولیه، دسته‌بندی جدیدی برای مقوله‌های زبانی تعریف شده است. دسته‌بندی ارائه‌شده، جز مواردی محدود، بسیار شبیه دسته‌بندی ارائه‌شده در روش گشتاری است.

دسته‌بندی پیشنهادی، مقوله‌های زبانی واژه‌های فارسی را در دو دسته مقوله‌های زبانی باز و مقوله‌های زبانی بسته قرار می‌دهد. مقوله‌های زبانی باز، فهرست بازی را تشکیل می‌دهند و در طول زمان ممکن است از آنها کاسته و یا به آنها اضافه شود. از طرف دیگر مقوله‌های زبانی بسته، فهرست بسته‌ای را تشکیل می‌دهد که تعداد آنها محدود و مشخص است و تنها برای نشان دادن نقش و یا رابطه دستوری خاصی به‌کار می‌رود، به عبارت دیگر استفاده از آنها جنبه نقشی دارد. در دسته‌بندی صورت‌گرفته، پاره‌ای از تعاریف سنتی در صرف فارسی، به دلیل ناکارآمدی در تحلیل رایانه‌ای، نادیده انگاشته شده است. این دسته‌بندی جدید در جدول های ۱ و ۲ نمایش داده شده است.

این دو گروه خود به دسته‌های کوچکتری تقسیم می‌شوند. این

دسته‌بندی‌ها با توجه به نحو و نقش متداول واژه و جایگاه آن در میان گروه‌ها و وابسته‌های پیشین و پسین در گروه‌های نحوی مختلف صورت گرفته است. علاوه بر اینکه دسته‌بندی فوق به شکلی صورت گرفته است که روابط منطقی بین گروه‌ها وجود داشته باشد. این روابط منطقی در قسمت بعد توضیح داده می‌شود. این روابط میان-گروهی، اساس کار موتور تحلیل صرفی را تشکیل می‌دهد. با مطالعه سطحی در دیگر زبان‌های هندی و اروپایی، دیده شده است که این روابط برای اکثر این زبان‌ها وجود دارد. تفاوت دیگر دسته‌بندی ارائه شده، قرارگیری مجموعه وندهای تصریفی مختلف در گروه های متفاوت است. این کار به خصوص در مورد زبان فارسی که در آن امکان دارد وندهای مقید جدا از پایه مربوط به آن در نوشته ظاهر شود، سبب سهولت در تحلیل صرفی و افزایش دقت خواهد شد.

در تحلیل نحوی، مقوله‌های زبانی ظاهر شده در زیرگروه مقوله‌های زبانی باز به‌همراه ویژگی‌های آنها به‌عنوان عناصر پایه‌ای دانش جهت تحلیل نحوی واژه‌ها استفاده می‌شود. مجموعه ویژگی‌های مقوله‌های زبانی از چگونگی ترکیب مقوله‌های زبانی باز و بسته فراهم می‌شود و به شکل زوج خصیصه-ارزش<sup>۱</sup> بازنمایی می‌شود.

جدول ۱- مقوله‌های زبانی باز

بن ماضی	بن مضارع	فعل
صفت	عدد	قید
اسم	ضمایر	شاخص
شبه جمله	حروف ربط	حروف اضافه

جدول ۲- مقوله‌های زبانی بسته

نشانه برتر برای قید	جز پیشین فعلی	نشانه‌های جمع
نشانه‌های فعلی	نشانه فعل دعایی	علامت نکره
الف ندا	نشانه‌های صفت مقایسه‌ای	ی واژه‌بست اضافی
نشانه‌های عدد ترتیبی	نشانه ماضی نقلی و بعید : ه	

### ۳- واژگان

در زبان‌شناسی لایه‌های متفاوتی برای بررسی زبان در نظر گرفته شده است. دانشی که در پردازش زبان به کار می‌رود، با توجه به این لایه عبارتند از: دانش واجی (دانشی درباره ساخت هجایی واژه در زبان)، دانش صرف (ساختواژی) (درمورد ساخت‌واژه در زبان)، دانش نحوی (درباره ساخت جمله زبان بر پایه واژه‌های آن زبان)،

دانش معنایی<sup>۲</sup> (درباره مفهوم و معنی پاره ای از گفتار)، و دانش کاربرد<sup>۳</sup> (درباره محل کاربرد واژه‌های زبان) و تحلیل کلام<sup>۴</sup> (درباره مفاهیم و ارتباطات با توجه به بافت زبانی)، تقسیم‌بندی می‌شود. هر یک از این لایه‌ها، از دانش ایجادشده در لایه‌ها بالاتر استفاده می‌کند.

واژگان<sup>۵</sup> مجموعه اطلاعات دانش واژگانی است که در پردازش زبان مورد استفاده قرار می‌گیرد. این دانش به واژه‌های زبان وابسته است. واژگان دانش اولیه و پایه در زبان‌شناسی رایانه‌ای را فراهم می‌آورد. در انجام این پژوهش مجموعه‌ای از واژگان، با توجه به دسته‌بندی ارائه شده در بخش ۲، تهیه شده است.

واژگان طراحی شده در پایین‌ترین سطح شامل یک پایگاه داده<sup>۶</sup> ساده است که هر یک از واژه‌های زبان در یک رکورد خاص جای گرفته اند. تراکنش‌های این پایگاه داده شامل اضافه نمودن واژه جدید به واژگان، حذف واژه و اصلاح یک واژه خاص است. برای بازنمایی واژه‌های زبان از یک ساختار خصیصه - ارزش استفاده شده است. در این روش هر واژه با یک مقوله زبانی علامت خورده است و هر مقوله زبانی ویژگی‌های خاص خود را دارد. واژه‌های هم‌نویسه<sup>۷</sup> که به یک صورت نگارش می‌شوند دارای مقوله‌های زبانی و یا معانی متفاوت می‌باشند و هر یک بصورت جداگانه در پایگاه داده ذخیره شده اند. همان‌طور که گفته شد دسته‌بندی صورت گرفته و ساختار مقوله‌های زبانی یا همان اجزای کلام سبب فراهم آمدن دانشی ضمنی بر روی واژگان شده است که از این دانش برای تحلیل لغوی استفاده می‌شود.

برای بازنمایی واژگان در سیستم پیشنهادی، از ساختمان داده درخت جستجوی سه‌تایی استفاده شده است. درخت جستجوی سه‌تایی اولین بار در [۲۱] معرفی شده است. ساختمان داده درخت جستجوی سه‌تایی به گونه ای است که رشته‌های اولیه مشترک میان داده‌های ذخیره شونده، تنها یک بار ذخیره می‌شود که نتیجه آن فشرده‌سازی مناسب در بازنمایی داده های متنی است. در این روش در هر گره از درخت، یک حرف ذخیره می‌شود. هر گره به سه گره چپ، وسط و راست اشاره می‌کند. اشاره‌گر سمت چپ به حرف با کد کوچکتر و اشاره‌گر سمت راست به حرف با کد بزرگتر اشاره می‌کند. اشاره‌گر وسط به کاراکتر بعدی در رشته ورودی اشاره می‌کند. پیمایش درخت توسط گره های سمت چپ و راست، سبب پیمایش در رشته ورودی نمی‌شود. در گره انتهایی هر واژه، مقوله واژه آن کلمه به همراه لیست خصیصه - ارزش ذخیره شده است. بازنمایی

<sup>2</sup> Semantic

<sup>3</sup> Pragmatic

<sup>4</sup> Discourse

<sup>5</sup> Lexicon

<sup>6</sup> Database

<sup>7</sup> Homograph

<sup>1</sup> Attribute-Value pairs

واژگان توسط این ساختار داده سبب فراهم آمدن ابزارهای قوی مانند جستجوی جزئی<sup>۱</sup> و همچنین سرعت بالا در جستجوی رشته جهت پیمایش رشته می‌شود.

#### ۴- روش پیشنهادی در طراحی تحلیلگر لغوی

هدف از تحلیل لغوی مشخص نمودن مقوله زبانی واژه و مجموعه ویژگی‌های صرفی آن است. دستیابی به این هدف نیازمند دانش وسیعی در محدوده صرف زبان می‌باشد. این دانش عموماً به کمک یکسری دانش رویه‌ای<sup>۲</sup> که به شکل قوانین نمود پیدا می‌کند و دانش تشریحی<sup>۳</sup> که در غالب واژگان ارائه می‌شود، فراهم می‌آید. با توجه به محدودیت‌های دانش رویه‌ای در ارائه دانش زبانی، روش‌های جدید همگی به سوی استفاده گسترده‌تر از دانش تشریحی که به شکل واژگان و پیکره‌های بسیار بزرگ نمود پیدا می‌کند، سوق داده شده‌اند.

تحلیل لغوی با صرف تصریفی در ارتباط است. از آنجا که صورت‌های صرفی هر واژه بسیار متنوع و زیاد است، در نتیجه امکان ذخیره آن در واژگان امکان‌پذیر نیست و کارایی لازم را به دنبال نخواهد داشت. خوشبختانه قوانین حاکم بر صرف تصریفی در زبان‌ها عمومی<sup>۴</sup> بوده و کمتر استثناپذیر است. به‌طور کلی، برای تحلیل یک واژه صرف شده، از یکسری قوانین که به‌طور متوالی بر واژه‌های زبان اعمال می‌شود، استفاده می‌گردد. اشکال این روش آن است که ترتیب و نوع اعمال قوانین انتخاب شده برای تحلیل واژه‌ها اهمیت دارد. اعمال نادرست قوانین و یا حتی ترتیب اشتباه سبب در بهترین حالت سیبی ناتوانی سیستم در تشخیص واژه و تحلیل آن می‌شود. علاوه بر این، جمع‌آوری قوانین و مشخص نمودن ترتیب مناسب جهت اعمال آنها، کاری پر زحمت و همراه با خطای فراوان است. در روش‌های سنتی زمانی تحلیل یک واژه به پایان می‌رسد که کلیه قوانین بر روی آن اعمال شده باشد. در بسیاری از موارد این کار سبب هدر رفتن زمان و منابع پردازشی می‌شود و حتی ممکن است در بعضی از موارد سبب ناپایداری سیستم به دلیل وجود قوانین زنجیره‌ای گردد. از سوی دیگر، در صورتی که شرط پایان یافتن تحلیل تنها یک جواب و یا به عبارت دیگر اولین جواب باشد، ممکن است قسمتی از اطلاعات از دست برود چراکه در بسیاری از موارد مانند واژه‌های هم‌نویسه یک واژه می‌تواند دارای چندین شکل صرفی و غیر صرفی باشد. به عبارت دیگر، روش‌های مبتنی بر قانون روش‌هایی فرضیه‌محوراند که از روش‌های بالا به پایین جهت تحلیل استفاده می‌کنند.

در اینجا و در روش پیشنهاد شده برای تحلیل‌گر لغوی، از نحوه

بازنمایی واژگان در درخت جستجوی سه‌تایی و دسته‌بندی ارائه شده برای مقوله‌های زبانی جهت انجام تحلیل لغوی استفاده می‌کند. درخت جستجوی سه‌تایی ابزار قدرتمندی را برای پیمایش رشته در اختیار کاربر قرار می‌دهد. روش پیشنهاد شده در این پژوهش از جستجوی جزئی و یک پایگاه دانش کوچک برای هدایت مسیر جستجو، جهت پردازش لغوی استفاده می‌شود. روش پیشنهاد شده، در مقابل روش‌های قبلی، روشی داده‌محور است. به عبارت دیگر پردازش لغوی واژه به مسئله پیمایش هدایت شده بر روی درخت جستجوی سه‌تایی تبدیل شده است.

ورودی سیستم تحلیلگر لغوی یک رشته از حروف زبان، در اینجا فارسی، است. در صورتی که رشته حروف وارد شده در سیستم، در زبان دارای معنی بوده و ریشه آن در واژگان ذخیره شده باشد با این تفاوت که شکل واژه ذخیره‌شده به دلیل تصریف تغییر یافته باشد، تحلیلگر لغوی، رشته حروف ورودی را پردازش خواهد نمود و نوع واژه و ویژگی‌های صرفی مرتبط با آن را مشخص می‌نماید.

با بررسی صورت گرفته بر روی واژه‌های زبان فارسی، به این نتیجه رسیدیم که جایگاه هر یک از اجزا از مقوله‌های زبانی در یک ترکیب که تشکیل یک واژه شده جدید را می‌دهند، ثابت است. می‌توان پس از یک واژه چندین جایگاه متفاوت در نظر گرفت. در هر یک از این جایگاه‌ها، تنها وندهای خاصی می‌توانند قرار گیرد. علاوه بر این که یکسری روابط درون گروهی بین مقوله‌های واژگانی که در درون یک ترکیب ظاهر می‌شوند وجود دارد.

ساختار درخت سه‌تایی امکان جستجوی جزئی را فراهم می‌آورد. به‌عنوان مثال اگر در واژگان واژه‌های 0,1,01,10 معرفی شده باشد آنگاه سیستم پس از مشاهده رشته 0110 متواند نتایج زیر را تولید کند: 0, 1, 1, 0, 1, 10, 0, 1, 0, 0, 0, 1, 01, 01, 10. با زیاد شدن تعداد واژه‌های زبان و افزایش طول رشته ورودی برای جستجوی جزئی، تعداد پیشنهادات سیستم و زمان پردازش به شدت افزایش می‌یابد. به عبارت دیگر، سیستم با یک فضای جستجوی بسیار بزرگ مواجه می‌شود. از میان پیشنهادات ارائه شده تنها یک یا چند مورد از آنها مورد قبول و مورد نظر است. افزایش فضای جستجو سبب افزایش زمان جستجو، ارائه جواب‌های اشتباه و در نهایت ناکارآمدی سیستم می‌شود. برای کاهش فضای جستجو از قوانین شهودی<sup>۵</sup> استفاده می‌شود. این قوانین به‌عنوان یک راهنما جهت محدود کردن فضای جستجو استفاده می‌شود. در سیستم پیشنهادی این کار به‌وسیله قوانین صرف انجام می‌شود. برای این منظور قوانین پیشنهادی می‌بایستی جامع و مانع باشند. مشکل بزرگ این مرحله، جمع‌آوری قوانین و حصول اطمینان از جامع و مانع بودن آنهاست. نحوه ارائه قوانین و اعمال آن، بر روی کارایی سیستم تاثیر به‌سزایی دارد. در سیستم پیشنهادی، این قوانین توسط دو ماتریس ارائه شده

<sup>1</sup> Partial Match

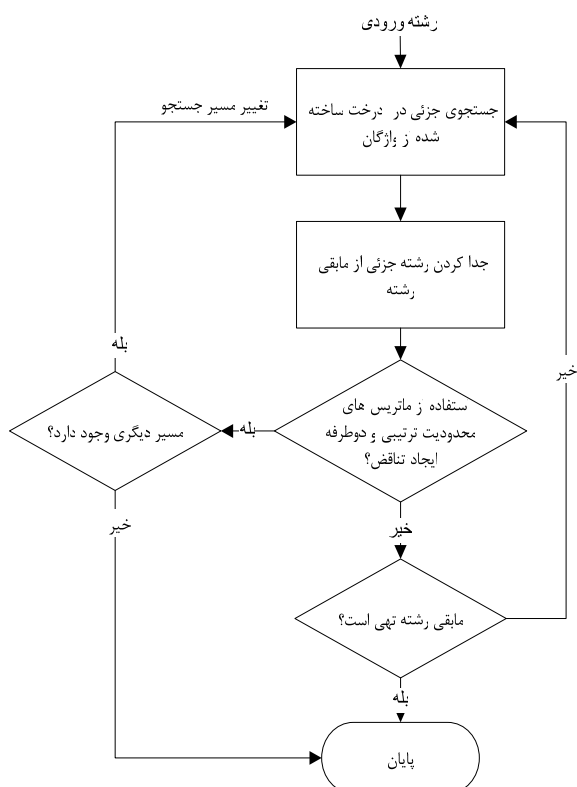
<sup>2</sup> Procedural Knowledge

<sup>3</sup> Declarative Knowledge

<sup>4</sup> General

<sup>5</sup> Heuristic

به‌عنوان مثال هیچگاه نمی‌توان شناسه را قبل از فعل آورد. این همان چیزی است که در سیستم‌های مبتنی بر قانون، به کمک ترتیب اعمال قوانین مورد توجه قرار می‌گرفت. در سیستم پیشنهادی چنین قوانینی در ماتریس محدودیت ترتیبی مدل‌سازی شده است. در اینجا نیز سطرها و ستون‌های ماتریس با مقوله‌های زبانی نام‌گذاری شده است. سطر و ستون متناظر با هر زوج مقوله زبانی نشان می‌دهد که آیا مقوله متناظر با یک سطر می‌تواند پس از مقوله زبانی مربوط به یک ستون بیاید یا خیر.



شکل ۱- بلوک دیاگرام الگوریتم پیشنهادی

روش کار بسیار ساده است. هنگامیکه یک رشته وارد سیستم می‌شود، رشته ورودی بر روی درخت جستجوی سه‌تایی ساخته شده از واژگان به شکل جزئی پیمایش می‌شود. هرگاه در این پیمایش جزئی، یک قسمت از واژه با یکی از واژه‌های واژگان سیستم نظیر شود، از مابقی رشته ورودی حذف می‌شود و به بقیه رشته‌های جزئی پیدا شده اضافه می‌شوند. در هر قدم پیش از اضافه شدن رشته جزئی به بقیه رشته‌های جزئی، از ماتریس‌های محدودیت دوطرفه و محدودیت ترتیبی برای اطمینان از عدم تناقض استفاده می‌شود. در صورت عدم وجود تناقض برای مابقی رشته در واژگان جستجو می‌شود تا هنگامی که به انتهای رشته ورودی برسد و یا با یک رشته پوچ مواجه شود. در صورت وجود تناقض، امکان پیمایش درخت توسط

که قوانین صرفی زبان به سادگی در آنها رمز شده است. این دو ماتریس، ماتریس‌های محدودیت دوطرفه<sup>۱</sup> و محدودیت ترتیبی<sup>۲</sup> نام گرفته‌اند. جدول شماره ۳ و ۴ قسمتی از این دو ماتریس را نمایش می‌دهد.

جدول ۳- قسمتی از ماتریس محدودیت دوطرفه

ی نکره	شناسه فعلی	جز پیشین فعلی	علامت جمع	اسم
1	0	0	1	اسم
1	0	0	1	صفت
1	0	0	1	اعداد اصلی
1	0	0	1	ممیزها
1	0	0	1	شاخص‌ها
0	0	0	1	ضمیر

همانطور که گفته شد، وجود یک مقوله زبانی در ساختمان یک کلمه سبب می‌شود تا یک محدودیت برای حضور مقوله‌های زبانی دیگر بوجود بیاید. به‌عنوان مثال در صورتی که در ترکیبی شناسه فعلی وجود داشته باشد آن‌گاه امکان حضور "ی" نکره در آن ترکیب نخواهد بود و بالعکس. مثالی دیگر برای زبان انگلیسی، حضور نشانه گذشته "ed" در انتهای یک فعل است که مانع از حضور نشانه استمرار "ing" در آن ترکیب می‌شود. ماتریس محدودیت دوطرفه همین واقعیت را نمایش می‌دهد. سطرها و ستون‌های این ماتریس با مقوله‌های زبانی نام‌گذاری شده‌اند. سطر و ستون متناظر با هر زوج مقوله زبانی نشان می‌دهد که آیا این دو مقوله می‌توانند با یکدیگر در یک ترکیب وجود داشته باشند یا خیر.

جدول ۴- قسمتی از ماتریس محدودیت دو طرفه

ی نکره	نشانه دعایی	شناسه فعلی	علامت جمع	ی نکره
1	0	0	0	علامت جمع
1	0	0	1	اعداد ترتیبی
0	1	0	0	شناسه فعلی
0	0	0	1	نکره
0	0	0	0	الف ندا
0	1	1	0	بن مضارع
0	1	1	0	بن ماضی

ترتیب ظهور هر یک از اجزای ترکیب در واژه‌سازی مهم است.

<sup>1</sup> Mutual restriction

<sup>2</sup> Order Restriction

مسیرهای دیگر چک می شود. در صورت عدم وجود مسیر دیگر جهت پیمایش، کار پایان می پذیرد. بدیهی است که هریک از رشته های جزیی با یک مقوله زبانی مشخص شده اند. همانطور که پیشتر گفته شد، از این مقوله ها و ماتریسهای محدودیت دو طرفه و محدودیت ترتیبی برای اطمینان از عدم وجود تناقض استفاده می شود. به عبارت دیگر، مسئله تحلیل صرفی به مسئله پیمایش درخت سه تایی تبدیل شده است. بلوک دیاگرام الگوریتم در شکل ۱ دیده می شود.

این الگوریتم به صورت بازگشتی پیاده سازی می شود و در حقیقت چک کردن قوانین صرفی در هنگام فراخوانی مجدد و با کمک ماتریسهای تعبیه شده صورت می گیرد. شکل زیر نمونه ای از خروجی سیستم برای دو واژه "کودکانشان" و "می زدنشان" را نمایش می دهد.

اسم	کودک	1	1	1
علامت جمع ان	ان	2	1	1
ضمیر ملکی	شان	3	1	1

جز پیشین فعلی	می	1	1	1
بن ماضی	زد	2	1	1
شناسه فعلی	ند	3	1	1
ضمیر مفعولی	شان	4	1	1

شکل ۳- نمونه خروجی سیستم برای دو واژه "کودکانشان" و "می زدنشان"

## ۶- ارزیابی عملی

سیستم پیشنهاد شده بر روی روزنامه شرق آنلاین در تاریخ ۸۴/۱/۲۹ آزمایش شد. پیکره مورد آزمایش مشتمل بر ۹۴۴۱۵ واژه کاربردی<sup>۱</sup> است. برای انجام این آزمایش از یک واژگان مشتمل بر ۲۲۰۰۰ واژه استفاده شد. در ۸۵٪ موارد واژه های پیکره عیناً در واژگان وجود داشت. ۱۲٪ واژه ها در پردازش لغوی تشخیص داده شد و تنها در ۳٪ موارد سیستم قادر به تشخیص واژه ورودی نبود. جدول شماره ۵ به طور خلاصه نتیجه آزمایش را با در نظر گرفتن واژه های تکراری نمایش می دهد. از سوی دیگر، واژه های استفاده شده در این پیکره مشتمل بر ۹۰۵۸ واژه نماینده<sup>۲</sup> است. ۵۹۵۰ واژه از واژه های نماینده عیناً در واژگان وجود داشت. تحلیل گر لغوی ۲۳۸۹ واژه از مابقی تعداد واژه های نماینده را تشخیص داد و ۷۱۹ واژه باقیمانده دیگر را تشخیص نداد. به عبارت دیگر، در ۷۱٪ موارد واژه عیناً در واژگان موجود بوده است. ۱۹٪ واژه ها توسط تحلیلگر لغوی تشخیص داده شده و در ۸٪ موارد، واژه وارد شده تشخیص داده نشده است. جدول شماره ۶ خلاصه این نتایج را نمایش می دهد.

جدول ۵- نتایج حاصل با احتساب واژه های تکراری

	تعداد	درصد
واژه تشخیص داده شده	80253	85
عیناً موجود در واژگان	11330	12
کلمات تشخیص داده نشده	2002	0.022
با شکل غیر صرفی	830	0.008
تعداد کل واژه ها (با احتساب واژه های تکراری)	94415	

جدول ۶- نتایج حاصل بدون احتساب واژه های تکراری

	تعداد	درصد
واژه تشخیص داده شده	6502	0.73
عیناً موجود در واژگان	1756	0.19
واژه تشخیص داده نشده	596	0.07
با شکل غیر صرفی	204	0.03
تعداد کل واژه ها (بدون احتساب واژه های تکراری)	8854	

میزان پوشش سیستم در تشخیص واژه های وارد شده با احتساب واژه های تکراری برابر ۹۷٪، و بدون احتساب واژه های تکراری در حدود ۹۳٪ بوده است.

از واژه های تشخیص داده نشده، با احتساب واژه های تکراری، ۸۳۰ مورد از واژه ها از ۲۸۳۲ مورد واژه تشخیص داده نشده شکل صرفی داشته اند. بدون احتساب واژه های تکراری، ۲۰۴ مورد از ۸۰۰ واژه تشخیص داده نشده، شکل صرفی داشته اند. فرمول شماره ۱ یک رابطه را برای ارزیابی دقت سیستم نمایش می دهد.

دقت سیستم = تعداد واژه های صرف شده در پیکره / فرمول (۱)  
تعداد واژه های صرف شده تشخیص داده شده توسط سیستم

مطابق این رابطه دقت سیستم با احتساب واژه های تکراری برابر ۹۳٪، و در صورت عدم احتساب واژه های تکراری برابر ۸۹٪ ارزیابی می شود. افزایش دقت به راحتی با اضافه نمودن ریشه واژه های تشخیص داده نشده امکان پذیر است.

## ۷- جمع بندی و کارهای آینده

در این مقاله روشی نوین جهت تحلیل لغوی واژه های زبان فارسی به کمک هماهنگی میان روش بازنمایی واژگان و الگوریتم پیمایش رشته با استفاده از ساختمان داده درختی ارائه شده است. در اینجا با استفاده از یک پایگاه دانش کوچک ارائه شده در غالب دو ماتریس که از قوانین صرف زبان فارسی و چگونگی هم نشینی مقوله های زبانی به دست آمده است، جهت هدایت پیمایش بر روی یالهای درخت

<sup>1</sup> Token

<sup>2</sup> Type

استفاده شده است. به عبارت دیگر، مسئله پردازش لغوی به مسئله پیمایش هدایت شده بر روی درخت جستجوی سه تایی تبدیل شده است. نتایج حاصل نمایانگر کارآمدی این روش برای تحلیل لغوی واژه‌های زبان فارسی است.

یکی از نقاط قوت روش ارائه شده، سادگی استفاده از روش پیشنهادی برای دیگر زبان‌های هندی و اروپایی است که در آنها تصریف واژه‌ها از طریق وندافزایی صورت می‌گیرد. مزیت دیگر این روش، سادگی جمع‌آوری دانش و ارائه آن می‌باشد بطور کلی ساختار پیشنهادی یک چارچوب کلی برای جمع‌آوری دانش صرفی زبان ارائه می‌نماید. سرعت بالا در پردازش، ارائه کلیه حالت‌های صرفی واژه به‌خصوص در مورد واژه‌های هم‌نویسه از دیگر مزیت‌های روش پیشنهاد شده است. در روش‌های قبلی برای تحلیل لغوی واژه‌های فارسی، عموماً قوانین کاربردی برای هر واژه براساس مقوله صرفی آن واژه مشخص می‌شود و بدین ترتیب تحلیل صرفی برای واژه‌ها به دو دسته تحلیل واژه‌های فعلی و غیر فعلی تقسیم بندی شده است. انتخاب قوانین براساس نوع وند ظاهر شده در واژه‌ها صورت می‌گردد و در صورتی که بخواهیم همگی صورت‌های ممکن برای تحلیل لغوی یک واژه را داشته باشیم، نتیجه عمل تحلیل در به‌کار بردن همگی قوانین و ترتیب مختلف آنها، سبب افزایش زمان تحلیل و افزایش حجم پردازش است. بر خلاف روشهای معمول، سیستم پیشنهادی نیاز به فرضیات اولیه جهت شروع به کار ندارد.

پژوهش صورت گرفته تنها در محدوده صرف تصریفی برای واژه‌های فارسی انجام شده است اما به راحتی امکان استفاده از چارچوب ارائه شده برای صرف اشتقاقی نیز وجود دارد. در این حالت کافی است تا سطر و ستون مربوط به وندهای اشتقاقی در زبان فارسی، به ماتریس‌های محدودیت ترتیبی و محدودیت دوطرفه اضافه شود.

یکی از مشکلات مهم که در طول کار با آن مواجه شدیم، عدم وجود استاندارد برای متون الکترونیکی فارسی است. این مشکل چه از نظر کد کارکترهای استفاده شده برای نمایش حروف و چه در نحوه تاپ و رعایت فاصله در مرز برونی واژه‌ها و نیم‌فاصله‌ها در مرز درونی واژه به چشم می‌خورد. شاید یک قدم بسیار مهم در این جهت، ارائه چارچوبی هماهنگ برای چگونگی بازنمایی متون الکترونیکی فارسی و ارائه استاندارد فراگیر باشد.

برای ادامه کار، استفاده از روش‌های یادگیری ماشین جهت تنظیم ماتریس‌ها و امکان رتبه بندی نتایج ممکن در نظر گرفته شده است. امید است که با این کار بتوان از واژه‌های هم‌نویسه در متون خاص، همزمان با تحلیل لغوی رفع ابهام صورت گیرد. همچنین، گسترش حجم واژگان جهت پوشش بیشتر از زبان مد نظر است. علاوه بر این که بسط و گسترش سیستم جهت تحلیل اشتقاقی واژه‌های فارسی، یکی از کارهای مد نظر در آینده می‌باشد.

## مراجع:

- [۱] ایران کلباسی، ساخت اشتقاقی واژه در فارسی امروز، پژوهشگاه علوم انسانی و مطالعات فرهنگی، تهران، ۱۳۸۰.
- [۲] کشانی، خسرو، اشتقاق پسوندی در زبان فارسی امروز، مرکز نشر دانشگاهی، تهران، ۱۳۷۱.
- [3] Kuhn, J. *Compounding and derivational morphology in a finite-state setting*, ACL, 2003.
- [4] Attia, M. A. *Developing a robust Arabic morphological transducer using finite state technology*, 8th Annual CLUK Research Colloquium, Manchester, 2005.
- [5] Vaillette, N. *Logical Specification of Finite-State Transductions for Natural Language Processing* PhD Thesis, Ohio University, 2004.
- [6] Radek, S. Pavel, S. *Automatic processing of Czech inflectional and derivative morphology*, FIMU-RS-2001-03, 2001.
- [7] Riazati, D. *Computational Analysis of Persian Morphology*, MS Thesis, Department of Computer Science, RMIT, 1997.
- [8] Yamashita, T. Matsumoto, Y. *Language independent morphological analysis*, In Proceedings of 6th Applied Natural Language Processing Conference, pp.232-238, 2000.
- [9] Bosch, A. Daelemans, W. , *Memory-Based Morphological Analysis*, ACL , 1999.
- [10] Snover M. G. Brent, M. R. *A probabilistic model for learning concatenative morphology*, NIPS, 2002.
- [11] Goldsmith J. A. *Unsupervised learning of the morphology of a natural language*, Computational Linguistics 27:2 pp. 153-198, 2000.
- [12] Sgarbas, K. N. Fakotakis, N. D. Kokkinakis, K. G. *A straightforward approach to morphological analysis and synthesis*, Proc. COMLEX , Greece, 2000.
- [13] Megerdoomian, K. *Unification-Based Persian Morphology*, Proceedings of the CICling, Mexico, 2000.
- [14] Evan A. L. *PC-KIMMO: A Two-Level Processor for Morphological Analysis*, Occasional Publications in Academic Computing No. 16. Dallas, TX: Summer Institute of Linguistics, 1990.
- [15] Megerdoomian, K. *Finite-State Morphological Analysis of Persian*, COLING, 2004.
- [۱۶] انوری، احمدی گیوی، دستور زبان فارسی ۲ (ویرایش دوم)، انتشارات فاطمی، تهران، ۱۳۸۲.
- [۱۷] طباطبائی، علاالدین، اسم و صفت مرکب در زبان فارسی، مرکز نشر دانشگاهی، تهران، ۱۳۸۲.
- [۱۸] ماهوتیان، شهرزاد، دستور زبان فارسی از دیدگاه رده شناسی، ترجمه سمائی، مهدی، نشر مرکز، تهران، ۱۳۸۲.
- [۱۹] دبیرمقدم، محمد، زبان شناسی نظری پیدایش و زایش دستور زایشی (ویراست دوم)، سازمان مطالعه و تدوین کتب علوم انسانی دانشگاهها(سمت)، تهران، ۱۳۸۳.
- [۲۰] مشکاةالدینی، مهدی، دستور زبان فارسی بر پایه نظریه گشتاری، انتشارات دانشگاه فردوسی مشهد، ۱۳۸۲.
- [21] Bentley J. L. Sedgewick, R. *Fast algorithms for sorting and searching strings*, Proceedings of the 8th Annual ACM-SIAM Symposium on Discrete Algorithms, 1997.