

# واژگان محاسباتی: ساختار مرکزی در سیستم‌های پردازش زبان طبیعی

احمد عبدالله زاده  
استادیار

مهرنوش شمس فرد  
دانشجوی دکترا

دانشکده مهندسی کامپیوتر، دانشگاه صنعتی امیرکبیر

## چکیده

امروزه مرکزیت واژگان در سیستم‌های پردازش زبان طبیعی مورد تأکید است و لذا طراحی و ساخت واژگان محاسباتی یکی از فعالیت‌های زیربنایی ساخت چنین سیستم‌هایی به شمار می‌آید. در این مقاله ابتدا به تعریف واژگان و اجزاء آن، بررسی عوامل افتراق واژگان‌ها و انواع حاصله می‌پردازیم. سپس بدنبال تشریح اصول نظری طراحی و ساخت واژگان‌ها، یک نمونه عملی را در یک زیر سیستم پردازش زبان فارسی ارائه خواهیم نمود. در این نمونه عملی، پارامترها و عوامل تأکید شده در مباحث نظری را مورد توجه قرار داده و مقدار دهی می‌نماییم.

## کلمات کلیدی

واژگان محاسباتی، پردازش زبان طبیعی، پایگاه دانش، هستان شناسی، یادگیری ماشینی

## Computational Lexicon: The Central Structure in Natural Language Processing Systems

M. Shamsfard  
Ph. D Student

A. Abdollahzadeh Barforoush  
Assistant Professor

Computer Engineering Department,  
Amir Kabir University of Technology

## Abstract

*Lexicons play a central role in natural language processing (NLP) systems, so design and implementing computational lexicons are of fundamental activities in development of such systems.*

*In this paper first, we define lexicon, introduce its parts, study the difference factors in lexicons and describe their various types. Then after discussing the theoretical aspects design and implementation of lexicons, we will study a practical instance in a Persian natural language processing system. In this system, we will consider the mentioned theoretical parameters and factors and initiate them practically.*

## Key words

*Computational Lexicon, Natural Language Processing, Knowledge Base, Ontology, Machine Learning.*

در این مقاله ابتدا پس از ارائه تعریفی از واژگان، اجزاء و همچنین وجوه افتراق انواع مختلف واژگان را بررسی خواهیم کرد. سپس مقایسه ای میان واژگان های محاسباتی و زبانی انجام داده و در نهایت به طرح ساختار و ویژگی های نمونه ای از واژگان های محاسباتی در یک سیستم پردازش زبان فارسی خواهیم پرداخت.

## ۲- تعریف واژگان

هر فرهنگ<sup>۴</sup> نگاشتی از «مجموعه نمادها» به «مجموعه تفاسیر» است. مجموعه نمادها می تواند شامل ترکیبات حروف الفباء ارقام، اصوات، تصاویر، اشکال و ... و یا دنباله ای از ترکیب آنها باشد و مجموعه تفاسیر بسته به نوع فرهنگ مورد نظر، معانی نمادها یا نمادهای دیگری متناظر با نمادهای داده شده را در بر می گیرد. به عبارت دیگر اگر فرض کنیم رابطه  $D$  یک رابطه فرهنگ سازی، مجموعه  $S$  مجموعه  $n$  عضوی نمادها و مجموعه  $I$  مجموعه  $m$  عضوی تفاسیر باشد، آنگاه خواهیم داشت:

$$D = \{(x,y) \mid x \in S, y \in I\}$$

$$S = \{x \mid x \text{ is a string of symbols;} \\ \text{symbol} \in A\}$$

$$I = \{y \mid y \text{ is a string of interpretations;} \\ \text{interpretation} \in B\}$$

$$A = \text{Alphabet} \mid \text{Digits} \mid \text{Phonemes} \mid \\ \text{Graphemes} \mid \dots$$

$$B = \text{Sememes} \mid \text{Alphabet} \mid \text{Digits} \mid \\ \text{Phonemes} \mid \text{Graphemes} \mid \dots$$

متداول ترین نوع فرهنگ، فرهنگ لغوی یا واژگان می باشد. در این نوع فرهنگ، مجموعه  $A$  مجموعه حروف الفبا و مجموعه  $B$  بسته به نوع کاربرد می تواند اجتماع مجموعه های واحدهای معنایی، الفبایی و صوتی باشد. واژگان در حقیقت انباره ای از اطلاعات در مورد زبان و ویژگی های اجتماعی آن است. یک واژگان اطلاعات لغوی مناسب و مرتبط (مانند طرز تلفظ و معنا) را برای واحدهای مناسب و مربوط لغوی<sup>۵</sup> (مانند لغات) ارائه می کند [۲]. معروف ترین شکل واحد لغوی کلمه است. اشکال دیگر آن می توانند قسمتی از کلمه یا ترکیبات کلمات باشند که در مورد آنها در بخش های بعد صحبت خواهیم کرد. هر واحد لغوی مجموعه ای از صور (تکواژها)<sup>۶</sup> است که مجموعه ای از محتواها (واحد های معنا)<sup>۷</sup> را پازنمایی می کند. این صور یا مجموعه ای

امروزه مرکزیت واژگان در سیستم های پردازش زبان مورد تأکید است. تقریباً بیش از ۱۰ سال است که محققین کار بر روی مسئله ورودی نامحدود را آغاز کرده اند، با این هدف که سیستم های پردازش زبان طبیعی بتوانند هر نوع جمله ای را دریافت کنند. به این منظور لازم است سیستم بتواند بطور وسیع واژگان و دستور زبان مورد نظر را ببوشاند. لذا طراحی ساختار واژگان، اکتساب دانش لغوی، روش های جستجو و ذخیره و بازیابی اطلاعات در واژگان و نحوه ارتباط واژگان با سایر اجزاء سیستم از مباحث مطرح در تحقیقات جهانی هستند.

در دهه ۸۰، منابع لغوی را دیکشنری ها (فرهنگ های لغات) تشکیل می دادند. فرهنگ های ماشین خوان<sup>۱</sup> (MRDs) فقط شکل تایپ شده فرهنگ های چاپی بودند که برخی اطلاعات فرهنگ نویسی اضافی مانند تغییرات قلم، دستورات تایپی و نمادهای خاص را در برداشتند. به مرور زمان نقش کامپیوترها در آماده سازی این فرهنگ ها بیشتر شد و جمع آوری داده ها، انتخاب و ساخت تعاریف، مرتب سازی مدخل ها و آزمون کامل و سازگار بودن اطلاعات وارده توسط خود ماشین صورت گرفت.

در دهه ۹۰، مهمترین منابع لغوی، پیکره های زبانی<sup>۲</sup> بوده اند. یک پیکره یک منبع اولیه<sup>۳</sup> است و حاوی مجموعه استراتژیکی از مواد زبان طبیعی (تمام متن یا نمونه هایی از متون با جملات بی ربط از یک یا چند زبان) می باشد که در شکل ماشین خوان ذخیره شده است. داده های این پیکره ها معمولاً پوشاننده حالات گفتاری و نوشتاری زبان است و محدوده وسیعی از زبان غیر فنی رایج (در زمان حال) را که به صورت عادی (و نه در شعر یا داستان های درام یا تخیلی) توسط بزرگسالان با لهجه استاندارد (و نه محلی)، استفاده می شود، دربر می گیرد.

امروزه واژگان های محاسباتی در گستره وسیعی از کاربردها بکار می آیند. از آن جمله می توان به سیستم های ویرایشگر، تصحیح املاء، تجزیه نحوی، تولید گرامر، ساخت و شناسایی گفتار، تولید و درک متن، طبقه بندی معنایی، رفع ابهام و ... اشاره نمود. برای ساخت یک واژگان محاسباتی دو راه اصلی وجود دارد ([۱]): راه اول بازسازی یک یا چند MRD و تبدیل آن به شکل مورد نیاز کاربرد است و راه دوم ساخت مستقیم واژگان از روی پیکره های زبانی می باشد. در بسیاری موارد از ترکیب این دو روش استفاده شده، ابتدا با استفاده از MRD ها واژگان اولیه ساخته و سپس بر اساس پیکره ها اطلاعات آن تصحیح و سازگار شده و یا تکمیل می گردد.

گرافیکی مثل حروف و یا مجموعه ای صوتی مثل صداها محقق می شوند. ارتباط بین صورت و معنا برای یک واحد لغوی یک ارتباط  $n$  به  $m$  است که  $n$  می تواند بزرگتر، کوچکتر و یا مساوی  $m$  باشد. به بیان دیگر:

$$L = \{(x, y) \mid x \text{ is a set of morphemes, } y \text{ is a set of sememes}\}$$

انواع دیگری از فرهنگ ها نیز وجود دارند که مورد توجه ما در این مقاله نیستند. مانند فرهنگ انتساب اشکال گرافیکی به حروف (مثلاً جهت بازشناسی متن)، انتساب کلمات به تصاویر (مثلاً برای استفاده کودکان)، انتساب اصوات به حروف (مثلاً برای تشخیص گفتار) و ...

### ۳- اجزاء واژگان

هر واژگان دارای دو جزء اصلی است:

- ساختار کلی

- رویه جستجو

ساختار کلی واژگان حاوی همه اطلاعاتی است که در واژگان جای گرفته است و خود از کنار هم قرار گرفتن اجزاء مستقلی بنام «مدخل»<sup>۸</sup> تشکیل شده است. یک مدخل، ساختاری است که اطلاعات مرتبط با یک واحد لغوی را در خود جای می دهد. این مدخل ها طوری مرتب شده اند که با یک رویه جستجو به سادگی قابل دستیابی اند. به بیان دیگر رویه جستجو امکان دستیابی سریع و آسان به مدخل های مرتب شده را فراهم می کند. نمونه ای از یک رویه جستجو می تواند ترتیب الفبایی باشد که بر اساس شکل نوشتاری واحد لغوی (و نه معنای آن) صورت می گیرد.

تصمیماتی نظیر تعیین اطلاعات مندرج در هر مدخل، نحوه کنار هم نشستن مدخل ها، چگونگی تخصیص واحدهای لغوی به مداخل و انتخاب رویه جستجوی مناسب بر عهده طراح واژگان است که بسته به نوع واژگان و کاربرد آن تعیین می نماید. اتخاذ تصمیمات متفاوت در هر زمینه می تواند منجر به ایجاد واژگان های مختلف گردد. نمونه هایی از این تفاوت ها در بخش بعد بررسی شده اند.

### ۴- انواع واژگان

واژگان ها بر اساس عوامل مختلفی طبقه بندی می شوند. از جمله این عوامل می توان به نوع و حدود اطلاعات مندرج در آنها و نحوه سازمان دهی این اطلاعات اشاره داشت. در دو بخش بعد عوامل طبقه بندی را بر اساس محتوا و ساختار واژگان خواهیم دید.

#### ۴-۱- تنوع در محتوا

محدوده و تنوع اطلاعات موجود در یک واژگان به عوامل مختلفی وابسته خواهد بود که الزاماً مستقل از یکدیگر نیستند. برخی از این عوامل در زیر آورده شده است.

**الف - درجه پوشا بودن واژگان:** این عامل نشان دهنده وسعت دانشی است که واژگان تحت پوشش قرار می دهد. یک واژگان از این جهت می تواند در دو سطح عمل کند: [۳]  
**الف - ۱- لغت نامه:** در چنین واژگانی هر لغت با کمک دانش زبانی موجود تعریف می شود. به عبارت دیگر لغات با کمک ارتباطاتشان با لغات دیگر مانند مترادف، تضاد و در برگیری معرفی می شوند. چنین واژگانی دانش جهان را که برای عملیات استنتاج و استدلال مورد نیاز است، در بر نمی گیرد.  
**الف - ۲- دائرة المعارف:** دایرة المعارف ها از جهت مجموعه تفاسیر یا لغت نامه ها متفاوتند. آنها بخش عظیمی از دانش جهان را در بر دارند. این دانش انجام استنتاج و استخراج دانش جدید را ممکن می سازد.

**ب- گروه مخاطبین واژگان:** واژگان بسته به این که برای چه گروه از مخاطبین و به چه منظور طراحی می شود دارای اطلاعات متفاوتی خواهد بود. دو نمونه از عوامل دسته بندی مخاطبین در زیر آمده است:

**ب- ۱- سن:** رده سنی مخاطب در انتخاب نوع توضیحات، رسانه مورد استفاده و ابزارهای بکار گرفته شده در تفسیر کلمات مؤثر است. به عنوان نمونه واژگان طراحی شده برای کودکان به زبان ساده تر، با توضیحات بسیط تر و اکثراً همراه با تصاویر و اشکال متنوع ارائه می شود. در چنین واژگانی مجموعه تفاسیر حاوی زیر مجموعه ای از اشکال و تصاویر است.

**ب- ۲- تخصص:** برخی واژگان ها تخصصی هستند و با در برگیری لغات و اصطلاحات فنی یک رشته، برای متخصصین آن رشته طراحی می شوند، در حالیکه برخی دیگر عمومی هستند و کلمات و لغات معمول زبان که افراد عادی از آنها استفاده فرهنگ ها هم در مجموعه نمادهای انتخابی و هم در مجموعه تفاسیر ارائه شده قابل مشاهده است.

**ج- تعداد زبان های تحت پوشش:** برخی واژگان ها تک زبانه هستند و برخی دیگر چند زبانه.

**ج- ۱- تک زبانه:** واژگان های تک زبانه لغات و اصطلاحات یک زبان را به همان زبان معنی کرده و یا در مورد آنها توضیحاتی ارائه می کنند. فرهنگ های لغات و دائرة المعارف ها نمونه هایی از چنین واژگان هایی هستند.

**ج- ۲- چند زبانه:** واژگان های چند زبانه لغات و عبارات مترادف در دو یا چند زبان را معرفی می نمایند. این فرهنگ ها بیشتر برای ترجمه بکار می روند. در چنین حالتی مجموعه نمادها از لغات زبان مبدا و مجموعه تفاسیر از لغات زبان

مقصد انتخاب می‌شوند. متداولترین این فرهنگ‌ها دیکشنری‌های دوزبانه هستند.

د- کاربرد واژگان: یکی دیگر از عوامل تنوع در واژگان نوع کاربرد آن است. اطلاعات مندرج در واژگان بسته به کاربرد و استفاده مورد نظر، متفاوت خواهد بود. به عنوان نمونه چند کاربرد در زیر آورده شده است:

د-۱- درک: واژگان مورد استفاده در درک متون لازم است از لحاظ معنایی و کاربردشناسی بسیار غنی باشد. این نوع واژگان بیشتر در سطح دائرةالمعارف عمل می‌کند و وظیفه آن نگاشت لغات و عبارات زبان به مفاهیم و معانی جهان واقعی است. در این حالت مجموعه تفاسیر، اجتماعی از مجموعه ویژگی‌های صرفی، نحوی و کاربردشناسی است.

د-۲- تولید: واژگان مورد نیاز از جهت ویژگی شبیه به مورد قبل است با این تفاوت که در جهت معکوس و به منظور نگاشت معانی و مفاهیم به لغات و عبارات زبان بکار می‌رود. به عبارت دیگر در این واژگان (نسبت به مورد د-۱)، جای مجموعه نمادها و مجموعه تفاسیر عوض شده است.

د-۳- ترجمه: واژگان مورد نیاز چند زبانه است و میزان پوشا بودن آن بستگی به سطح ترجمه مورد نظر دارد. در این حالت مجموعه نمادها و تفاسیر از الفبای دو یا چند زبان متفاوت گرفته می‌شوند.

ه- کاربرد واژگان: عامل مؤثر دیگر در تعیین محتوای واژگان، کاربرد آن است. اطلاعات عمومی و قبلی کاربر در هنگام استفاده از واژگان، تعیین‌کننده سطح اطلاعاتی است که لازم و کافی است در واژگان جای گیرد.

ه-۱- انسان: در صورتی که کاربر انسان باشد، سطح وسیعی از دانش عرفی و عمومی را داراست. لذا می‌توان با استفاده از اصل کم‌کوشی، از درج مجدد اطلاعاتی که فرض می‌شود مخاطب دارد خودداری نمود. در این حالت انسان با دانش عمومی خود دانش مندرج در واژگان را تفسیر و خلاءهای اطلاعاتی آن را بسهولت پر می‌کند. در چنین واژگانی بخشی از مجموعه تفاسیر و همچنین زیر مجموعه‌ای از مجموعه نمادها، دانسته فرض شده و از واژگان حذف می‌شوند.

ه-۲- ماشین: اگر کاربرد واژگان ماشین یا سیستم کامپیوتری باشد، لازم است دانش عرفی و عمومی مورد نیاز نیز در واژگان تعبیه شود. لذا حجم و محتوای دانش مندرج در واژگان محاسباتی که کامپیوتر از آن استفاده می‌کند با واژگان زبان شناسی مورد استفاده انسان‌ها متفاوت است. در واژگان محاسباتی، مجموعه تفاسیر و نمادها بزرگتر از واژگان زبان شناسی بوده، از جهت صورت بندی نیز بدلیل محدودیت‌های قابلیت خوانایی توسط ماشین، متفاوت خواهند بود.

## ۴-۲. تنوع در ساختار و نحوه انجام عملیات

سازمان دهی اجزاء یک واژگان بر اساس پارامترهای مختلف و به صور متفاوتی انجام می‌شود. این تفاوت در ساختار منجر به ایجاد تنوع در واژگان‌ها خواهد شد. برای سازمان دهی واژگان، طراحی ساختار و الگوریتم‌های انجام عملیات اصلی در واژگان، لازم است طراح به سؤالاتی پاسخ دهد. مجموعه پاسخ‌های طراح به این سؤالات، ساختار و اعمال لازم برای واژگان را تعیین خواهد کرد. نمونه‌هایی از موارد سؤال و برخی پاسخ‌های نوعی به آنها در زیر آورده شده است:

### الف- نحوه برخورد با کلماتی با چند معنی

در هر زبانی تعدادی از کلمات دارای بیش از یک معنی هستند. نحوه سازمان دهی این معانی برای کلمه یکی از مسائلی است که طراح باید به آن پاسخ دهد. به این منظور دو امکان بررسی می‌شود:

الف-۱- آوردن همه معانی تحت یک مدخل: در این روش برای هر کلمه با هر تعداد معنی تنها یک مدخل در نظر گرفته می‌شود و تمامی معانی کلمه در همان مدخل جای می‌گیرند. در این حالت واژگان تابعی از مجموعه نمادها به مجموعه تفاسیر است. به عنوان مثال در چنین واژگانی برای کلمه «شیر» تنها یک مدخل وجود دارد که در مقابل آن معانی چون: حیوانی درنده، وسیله خروج آب یا مایعات از لوله و مایع خوراکی مترشحه از سینه پستانداران، درج می‌شود.

الف-۲- استفاده از مدخل‌های مختلف: در این روش معانی مختلف کلمه بر اساس معیارهای مختلف در مدخل‌های جداگانه جای داده می‌شوند. این معیارها می‌توانند داشتن ریشه کلمه، منشأ تاریخی، ویژگی‌های نحوی یا کاربردشناسی متفاوت باشند. در این حالت واژگان دیگر یک تابع از مجموعه نمادها به مجموعه تفاسیر نیست بلکه رابطه‌ای است که عضو اول آن لزوماً یکتا نیست. در این واژگان برای همان کلمه «شیر» سه مدخل جداگانه در نظر گرفته می‌شود که هر یک حاوی یکی از معانی کلمه است.

الف-۳- حالت ترکیبی: این حالت که معمولاً متداول‌تر از حالات قبل است ترکیبی از حالات الف-۱ و الف-۲ است. در این واژگان معیارهای تفکیکی محدود شده و تنها به ازای برخی معیارها مثلاً منشأ تاریخی کلمه مدخل‌ها مجزا و سایر تفاوت‌ها تحت یک مدخل بیان می‌شوند.

### ب- نحوه نگاشت مدخل‌ها به واحدهای لغوی

ب-۱- یک به یک: در این حالت هر واحد لغوی با هر تعداد معنی یک و فقط یک مدخل را به خود اختصاص می‌دهد (مشابه حالت الف-۱) و هر مدخل نیز تنها به یک واحد لغوی تخصیص می‌یابد.

ب- ۲- یک به چند: در این حالت یک مدخل به چند واحد لغوی تخصیص می‌یابد. در چنین شرایطی یک واحد لغوی مدخل اصلی محسوب شده و چندین واحد دیگر تحت این مدخل به عنوان زیر مدخل یا وابسته‌های آن قرار می‌گیرند.

ب- ۳- چند به یک: در این حالت چند مدخل برای یک واحد لغوی صرف می‌شوند. این حالت وقتی اتفاق می‌افتد که یک واحد لغوی با چند معنی داریم که برای هر معنی یک مدخل در نظر گرفته ایم (مشابه حالات الف- ۲ و الف- ۳).

### ج- نحوه تعیین واحدهای لغوی

ج- ۱- کلمات: ساده‌ترین و بدیهی‌ترین شکل واحدهای لغوی کلمات هستند. مانند خانه، آوردن و ...

ج- ۲- زیر کلمات: زیر کلمات تکواژه‌هایی هستند که به ساخت تصریف‌ها یا مشتقات کمک می‌کنند مانند پیشوندها، پسوندها و میانوندها.

ج- ۳- زیر کلمات: ترکیباتی که از بیش از یک کلمه درست شده‌اند نیز می‌توانند واحد لغوی باشند. این ترکیبات در گروه اصلی جای دارند:

- واحدهای ترکیبی و آغازگرا: هر دو نوع این واحدها اجزائی از چند کلمه دارند با این تفاوت که در واحدهای ترکیبی این اجزاء از هر کجای کلمه می‌توانند انتخاب شوند ولی در واحدهای آغازگرا این اجزاء از حروف اول کلمات برگرفته می‌شوند. به عنوان مثال می‌توان به واحد «مشاصد» (ترکیبی از حروف کلمات مشاورین صنایع دفاع) به عنوان یک واحد ترکیبی و به واحد «نرزا» (ترکیبی از حروف آغازین کلمات نیروی زمینی ارتش جمهوری اسلامی ایران) به عنوان یک واحد آغازگرا اشاره نمود.

- اصطلاحات و عبارات مرکب: این واحدها ترکیبی از چند کلمه کامل هستند. راه تشخیص این واحدها از دنباله‌ای از کلمات متوالی این است که در این حالت معنی کل عبارت با مجموع معانی کلمات تشکیل دهنده آن متفاوت است. مثل واحد «سر دواندن» یا اصطلاحات و ضرب المثل‌ها.

### ج- ۴- موارد خاص شامل

اسامی خاص: یک گروه مهم از مدخل‌ها را اسامی خاص تشکیل می‌دهند. این مدخل‌ها ممکن است به اسامی اشخاص (مثل ارسطو یا فردوسی)، مکان‌ها (مانند حافظیه یا خزر) یا اصطلاحات عملی و فنی (مانند پنی سیلین یا رایانه) تخصیص یابند. برخی واژگان‌ها برای محدود نمودن مجموعه لغات خود این مدخل‌ها را حذف می‌نمایند.

صورت‌های روزمره: این عبارات واحدهای غیر لغوی هستند که کاربرد عملی دارند یا ترجمه واحد لغوی در زبان دیگر می‌باشند (مانند ترجمه استعارات).

### د- یگانگی ساختار واژگان

د- ۱- تک ساختاری: در این روش همه واحدها در یک ساختار ریخته می‌شوند. شکل واژگان برای همه مدخل‌ها یکسان بوده، واژگان یکدست و بدون ضمایم ارائه می‌شود.

د- ۲- چند ساختاری: در روش چند ساختاری در کنار ساختار اصلی واژگان، از ضمایم برای معرفی واحدهای خاص مانند اسامی خاص و مخفف‌ها استفاده می‌گردد. ساختار این ضمایم و نوع اطلاعات مندرج در آنها می‌تواند به کلی با ساختار اصلی واژگان متفاوت باشد.

### ه- تعیین مدخل‌های اصلی و زیر مدخل‌ها

متداول‌ترین ترتیب نزولی بر اساس احتمال مدخل بودن و ترتیب صعودی بر اساس احتمال زیر مدخل بودن در انگلیسی به شکل زیر است: تکواژها (مثل home)، واحدهای ترکیبی (مثل smog)، واحدهای آغازگرا (مثل VIP)، ترکیبات اسمی با نوشتار بسته (بدون فاصله) (مثل Black-bird)، ترکیبات اسمی با نوشتار باز (با فاصله) (مثل night owl)، ترکیبات فعلی (افعال مرکب) (مثل give up)، اصطلاحات غیر فعلی (مثل at all)، اصطلاحات فعلی (مثل kick the bucket) و مشتقات با معنای شفاف (که معنای آنها از معنای اجزایشان قابل درک است) شامل: ترکیبات پسوند (زیر قسمت اصلی کلمه) (مثل lameness تحت مدخل lame)، و پیشوندی (به ترتیب الفبایی زیر پیشوند) یا به ترتیب غیر الفبایی زیر قسمت اصلی کلمه با ارجاع به محل اصلی الفبایی آنها) (مثل pre-war تحت مدخل pre-).

### و- ترتیب واحدهای با تشابهات صوری

در هر زبانی کلماتی وجود دارند که از لحاظ شکل نوشتاری بهم شبیه هستند ولی معانی متفاوتی دارند. این تشابهات صوری بر دو گونه‌اند: دسته اول کلمات با چند معنی هستند. در این حالت نه تنها شکل نوشتاری بلکه تلفظ کلمه نیز در معانی مختلف ثابت می‌ماند. (مانند کلمه شیر با سه معنی متفاوت) اما دسته دوم کلمات با شکل نوشتاری یکسان و تلفظ متفاوت هستند (مانند کلمه حکم با تلفظ هـ یا hokm, hakam). هر دو دسته در گروه تشابهات صوری قرار می‌گیرند. نکته قابل توجه طراح در مورد این تشابهات تعیین ترتیب وقوع این واحدها در واژگانی است که بر اساس شکل صوری کلمات مرتب شده است. به این منظور معیارهای مرتب سازی زیر پیشنهاد می‌گردد:

و- ۱- به ترتیب تاریخی: در این روش واحدهای مشابه بر اساس منشاء تاریخیشان (معمولاً از قدیم به جدید) مرتب می‌شوند.

و- ۲- به ترتیب فراوانی: در این روش واحدهای پرمصرف تر پیش از واحدهای مشابه کم مصرف آورده می‌شوند. به عبارت دیگر بر اساس مطالعات آماری، واحدها بر مبنای

بیشترین رخداد معانی مرتب می شوند.

و- ۳- به ترتیب الفبایی اجزاء کلامشان: راه دیگر مرتب سازی واحدهای مشابه در نظر گرفتن ترتیب الفبایی طبقه نحوی آنهاست. مثلاً اسم و صفت پیش از فعل و قید پس از آن ظاهر می شود.

و- ۴- استفاده از الگوریتم های خاص مثلاً

- قائل شدن ترتیب بین حروف کوچک و حروف بزرگ

- قائل شدن ترتیب بین واحدهای بسته (بی فاصله) و واحدهای باز (فاصله دار)

- قائل شدن ترتیب بین حروف خاص مثل خط فاصله و آپاستروف

ز- اهم اطلاعات ذخیره شده برای هر مدخل (واحد لغوی) یکی از مهمترین مسائلی که طراح باید در مورد آن تصمیم بگیرد، انتخاب اطلاعاتی است که لازم است برای هر واحد لغوی ذخیره شود. در زیر برخی از این اطلاعات آورده شده است. واژگان های مختلف بسته به نوع و کاربردشان ممکن است برخی یا همه این اطلاعات را در خود جای دهند.

ز- ۱- شکل واحد (با گونه های نوشتاری مختلف): همان رسم الخط واحد لغوی است که معمولاً به عنوان اندیس مدخل هم استفاده می شود. مانند: برد، کتاب، ...

ز- ۲- اطلاعات آوایی (تلفظ و الگوهای تأکید): مانند: (برد: bord)، (کتاب: keta'b)

ز- ۳- طبقه بندی و زیر طبقه بندی های نحوی شامل: اجزاء کلام (اسم، فعل، صفت و غیره): مثال: (برد: فعل)، (کتاب: اسم)

- اطلاعات دستوری بیشتر (شمارش پذیری، تعدی و غیره): مثال: (برد: ماضی، متعدی)، (کتاب: مفرد، شمارش پذیر) ز- ۴- گونه های تصریفی (جمع، ماضی، اشکال بی قاعده و ...): مثال: (برد: مضارع: برد: barad)، (کتاب: جمع: کتب: kotob)

ز- ۵- گونه های اشتقاقی و ترکیبی: مثال: (برد: برد و باخت: نتیجه مسابقه)، ...، (کتاب: مکتوب: نوشته شده)، (مکتب: ۴۱- محل نوشتن، ۲- نوع نگرش خاص)، (سر کتاب باز کردن: دعا نویسی)...

ز- ۶- معرف معنی (ها) با ارائه تعاریف یا ارجاع به مترادف ها: مثال: (برد: ۱- انتقال داد، ۲- پیروز شد)، (کتاب: مجموعه چاپی از اوراق نوشته شده)

ز- ۷- مثال هایی از طرق مختلف استفاده از کلمه: (برد: ۱- علی کتاب را به مدرسه برد. ۲- علی مسابقه را برد. علی شرط را برد.)، (کتاب: کتاب را خواندم.)

ز- ۸- موارد استفاده یا استفاده خاص: مثال: (برد: از میان برد: او آبروی مرا برد)

ز- ۹- اطلاعات در مورد منشأ واحد: در مورد ریشه کلمه و پیدایش و تغییر شکل آن در طول زمان تا رسیدن به شکل فعلی صحبت می کند.

ز- ۱۰- اطلاعات جانشینی<sup>۹</sup> مانند مترادف، متضاد، وارون، فراگیر و ...: مثال: (برد: مترادف: انتقال داد)، (فراگیر: انتقال فیزیکی)، (وارون: آورد)، (کتاب: فراگیر: شیء فیزیکی - بی جان)

ز- ۱۱- اطلاعات هم نشینی<sup>۱۰</sup> مانند عبارات و ترکیباتی که این کلمه معمولاً در آنها واقع می شود و یا محدودیت های گزینشی آن: مثال: (برد: چیزی یا کسی را بردن: فاعل: متحرک، مفعول: قابل حمل)، (کتاب: کتاب خواندن، کتاب نوشتن: فاعل: انسان، باسواد)

ز- ۱۲- اطلاعات تشابهی<sup>۱۱</sup> (قیاس پذیر) حاوی اطلاعاتی در مورد حوزه لغوی که این واحد را شامل می شود (مانند استفاده از see also)

ز- ۱۳- اطلاعات نظام مندی<sup>۱۲</sup> زبان مانند تعلق واحد به هسته زبان یا قرصی بودن آن، تعلق به زمان یا مکان و یا قلمروی تخصصی خاص، و یا رسمی یا غیر رسمی بودن واحد: مثال: (برد: متعلق به هسته زبان، عام)

ح- تکنیک های انجام توضیحات:

ح- ۱- نمایش: با جدول، شکل یا دیاگرام

ح- ۲- مثال آوری: ارائه مثال های مختلف از کاربرد واحد لغوی و نحوه ترکیب آن با واحدهای دیگر

ح- ۳- بسط: این نوع توضیح برای مخفف ها یا واحدهای ترکیبی و آغازگرا بیشتر بکار می رود و شکل کامل عبارتی را که واحد لغوی از آن گرفته شده بیان می کند. گاهی بیان شکل کامل عبارت اطلاعات کافی در اختیار کاربر می گذارد و نیازی به توضیح بیشتر نیست مانند توضیح واحد نزاچا به شکل نیروی زمینی ارتش جمهوری اسلامی ایران. اما بعضی مواقع بسط واحد به تنهایی اطلاع زیادی در مورد آن به ما نمی دهد مثلاً بسط واحد ناتو (NATO) به شکل "North Atlantic Treaty Organization" است که اطلاعی در مورد این سازمان، اعضاء عملکرد و اهداف آن نمی دهد.

ح- ۴- بحث: توضیح ساده استفاده روزانه از این واحد

ح- ۵- تعریف: شکستن واحد به اجزاء سازنده و ترکیب مجدد آنها طوری که محتوای واحد ساخته شده، معنی یا تعریف واحد مورد نظر باشد. تعریف ممکن است به صورت تعریف با مترادف، تعریف تحلیلی و یا تعریف فرمولی باشد.

## ۵- مقایسه واژگان محاسباتی و واژگان زبان

واژگان محاسباتی و واژگان زبان از جنبه های مختلف با هم متفاوتند. مهمترین و واضح ترین تفاوت در انتظاری است

که از آنها داریم. در واژگان زبان (اعم از شکل چاپی و یا تایپ شده آن) کاربر انسان است. لذا اطلاعات مندرج به زبان طبیعی و با صورت بندی قابل فهم برای انسان در واژگان ذخیره می شوند. در مقابل، کاربر اصلی واژگان محاسباتی، ماشین است. لذا اطلاعات با فرمت قابل خواندن توسط ماشین و عموماً به یک زبان بازنمایی میانی نگهداری می شود. در این سیستم ها در صورت نیاز، واسطه های کاربری تعبیه می شود که اطلاعات مورد نیاز کاربر انسانی را تغییر شکل داده به صورت قابل فهم توسط انسان به او ارائه نماید.

به دلیل اختلاف فوق واژگان های محاسباتی و زبانی هم از جهت محتوایی و هم از جهت ساختاری با هم متفاوتند. اختلاف محتوایی واژگان محاسباتی و زبان شناسی ناشی از اختلاف در کاربر آنهاست. درجایی که کاربر واژگان انسان است، حجم زیادی از دانش عرفی و عمومی بشر، با فرض علم خواننده نسبت به آن صراحتاً آورده نمی شود. این بخش از دانش که برای انسان دانسته فرض شده، برای ماشین ناشناخته است. لذا در واژگان های محاسباتی لازم است صراحتاً بیان گردد. این تفاوت بیشتر در دانش های معناشناسی و کاربردشناسی نمود دارد.

اما اختلاف ساختاری واژگان های زبان شناسی و محاسباتی از چند جهت قابل بررسی است:

اولاً واژگان محاسباتی بدلیل دارا بودن حجم بیشتری از دانش، نیاز به ساختار بزرگ تر و پیچیده تری جهت ذخیره سازی این دانش دارد.

ثانیاً در واژگان محاسباتی برای دستیابی سریع و چند جهته به اطلاعات، اتصالات و ارتباطات بیشتری (نسبت به واژگان زبان شناسی) تعبیه شده اند. به عنوان نمونه در یک واژگان غیر محاسباتی که از سلسله مراتب ارث بری برای نمایش معنا استفاده می کند، معمولاً هر گره به گره پدرش اشاره دارد ولی به فرزندان و برادرانش اشاره نمی کند. مثلاً در تعریف کلمه «درخت» گفته می شود که «درخت» یک نوع «گیاه» است و در تعریف کلمه «گیاه» ارتباط فرزندی با گره «موجود زنده» قید می گردد. اما در تعریف گره «درخت» به فرزندان دیگر گره پدر (برادران) مانند «گل» اشاره نمی شود. همچنین ارتباط معکوس میان گره پدر و فرزند برقرار نیست. مثلاً از همان گره «درخت» راهی برای رسیدن به گره های «راش» و «افرا» وجود ندارد. در حقیقت این اطلاعات مفقوده اطلاعات ساختاری هستند. چرا که به هر حال از روی واژگان غیر محاسباتی نیز می شود آنها را استخراج کرد ولی کار بسیار مشکل و زمان بری خواهد بود. در حالیکه با افزایش اطلاعات ساختاری و ارتباطات میان اجزاء واژگان، سرعت دستیابی به این اطلاعات را افزایش داده ایم. نمونه دیگر این

ارتباطات اضافه می تواند تعیین معنی ای از گره پدر باشد که گره جاری فرزند آن است. در بسیاری موارد گره پدر دارای بیش از یک معنی است. در واژگان های غیر محاسباتی، در چنین حالتی گره های فرزند تنها به گره پدر خود به شکل عام اشاره دارند و معنی مورد نظر خود از گره پدر را تعیین نمی کنند. در بسیاری از واژگان های محاسباتی (و از جمله در وردنت [۵]) برای افزایش سرعت استخراج اطلاعات چنین ارتباطاتی بوجود می آیند.

ثالثاً بدلیل استفاده ماشین از واژگان محاسباتی، لازم است این واژگان به صورت کاملاً ساخت یافته و با ساختار از پیش تعیین شده ارائه گردد. به همین دلیل واژگان های مختلف، جهت تطابق ساختارهای لازم برای نمایش همه طبقات نحوی و معنایی زبان، تدابیر متفاوتی اندیشیده اند. مثلاً وردنت از واژگان چندساختاری بهره گرفته، ساختارهای متفاوتی برای بازنمایی اطلاعات در مورد اعضاء طبقات نحوی مختلف در نظر می گیرد. در برخی واژگان های محاسباتی دیگر، تک ساختاری ترجیح داده شده و ساختمان داده های مورد نیاز به شکل پویا و قابل انطباق تعریف می شوند. این در حالی است که اکثر واژگان های غیر محاسباتی تک ساختاری هستند و از رویه ای مسطح برای معرفی همه واحدهای لغوی استفاده می کنند (یا حداقل اگر چندساختاری باشند، معیار تقسیم بندی کاملاً متفاوت با واژگان های محاسباتی دارند). در چنین واژگانی اصراری برای انطباق ساختاری و یکسان سازی اجزاء تعریف همه واحدهای لغوی وجود ندارد.

تفاوت دیگری که هم از جهت ساختاری و هم از جهت محتوایی قابل بررسی است، تفاوت در نحوه ارائه معانی در دو نوع واژگان است. در واژگان زبان شناسی، معمولاً هر واحد لغوی با معرفی مترادف ها، بیان توضیحات و یا ارائه تصاویر معنی می شود. استفاده از کلمات ساده تر برای معرفی کلمات و عبارات دیگر بدلیل احتمال زیاد شناخت کاربر از کلمات ساده، معمول است. در چنین واژگانی گاهی بدلیل استفاده از کلمات برای تعریف کلمات دیگر، با حلقه مواجه می شویم. اما در واژگان های محاسباتی نحوه معرفی معنی متفاوت است. در بعضی سیستم هایی که از واژگان بهره می گیرند (مانند مایکروکاسموس [۶])، برای نمایش معنی از ساختار دیگری به نام هستان شناسی استفاده می شود. هستان شناسی دانش مفاهیم جهان را درون خود سازمان دهی می نماید. معنی نمودن هر واحد لغوی در واژگان با معرفی جایگاه آن در سلسله مراتب هستان شناسی انجام می پذیرد. به این ترتیب کلمات مترادف آنهایی هستند که به مفاهیم و معانی یکسان در هستان شناسی ارجاع داشته

می‌شود. به عبارت دیگر واژگان، ساختاری با محتوای پویاست. برای طراحی واژگان لازم است نوع واژگان، اطلاعات مندرج در آن، ساختار مورد نیاز و نحوه دستیابی، ذخیره سازی، بازیابی و بهنگام سازی این اطلاعات مشخص گردد. در ادامه این بخش به بررسی مشخصات مذکور در طراحی واژگان در سیستم «هستی» می‌پردازیم.

#### الف - نوع واژگان

- درجه پوشا بودن واژگان این سیستم خود به تنهایی به صورت یک لغت نامه عمل می‌کند. اما از آنجا که در ارتباط با هستان شناسی قرار می‌گیرد و از طریق مدخل‌های آن می‌توان به دانش گسترده مفاهیم و ارتباطات جهان دسترسی داشت، در سطح دائرة المعارف مورد استفاده قرار می‌گیرد. - گروه مخاطبین: مخاطب این واژگان در حقیقت سیستمی است که قرار است از متون در حد کتاب فارسی کلاس اول دبستان، اطلاعاتی را استخراج کند. لذا در برگزیده واژه‌هایی در حد همین کتاب است و واژه‌ها و عبارات تخصصی در آن جای ندارند.

- تعداد زبان‌های تحت پوشش: واژگان مورد نظر در هستی تک‌زبانه است و فقط اطلاعات در مورد واژه‌های زبان فارسی را در بردارد.

- کاربرد واژگان: این واژگان فعلاً جهت درک متون به کار می‌رود اگر چه می‌توان در مراحل بعدی پروژه و یا در گسترش‌های آتی آن به مسئله تولید نیز توجه نمود.

- کاربرد واژگان: این واژگان یک واژگان محاسباتی است و کاربرد آن ماشین می‌باشد.

- نحوه برخورد با کلماتی با چند معنی: در این واژگان، برای حفظ نمود<sup>15</sup>ها (مفهوم‌ها)ی مختلف یک کلمه از یک مدخل با فلیدهای مختلف استفاده می‌شود. این تصمیم براساس ساختار انتخابی برای واژگان و امکانات زبان برنامه نویسی منتخب [۱۰] اتخاذ شده است. کلمات یکسان نه تنها ممکن است معانی مختلفی داشته باشند (مثل سه معنی کلمه شیر)، بلکه ممکن است دارای ویژگی‌های لغوی و نحوی متفاوتی نیز باشند (مثل شکست: فعل ماضی سوم شخص مفرد و

باشند. با این روش اگر چه مشکل معرفی مفاهیم در هستان شناسی اضافه می‌شود، ولی امکان بروز حلقه از میان رفته کار معرفی معنا در واژگان تسهیل می‌گردد. در برخی واژگان‌های دیگر نیز که از هستان شناسی بهره نمی‌گیرند (مانند وردنت) معرفی معنی واحد لغوی با بیان ارتباطات مختلف آن با واحدهای دیگر انجام می‌گیرد. وسعت این ارتباطات منجر به حذف حلقه و رفع ابهام در واژگان می‌گردد.

### ۶- طرح یک نمونه عملی: واژگان در پروژه هستی

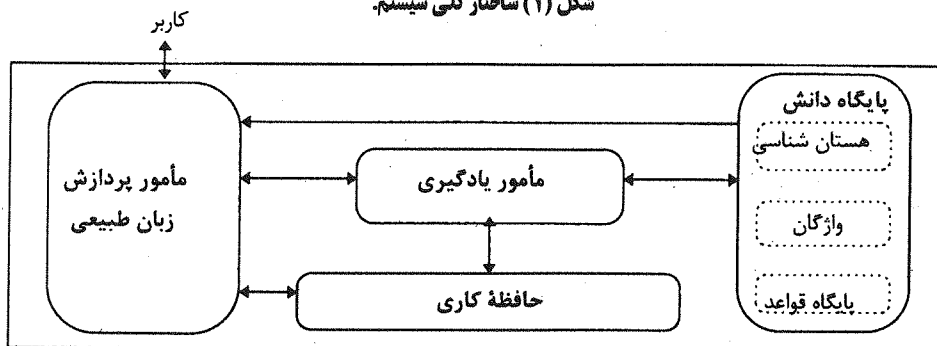
#### ۱-۶- معرفی هستی

هستی [۸ و ۷] سیستم ساخت خودکار هستان شناسی<sup>13</sup> بر اساس پردازش زبان فارسی است. این سیستم از مدل یک هسته یادگیر و تکاملی برای ساخت پایگاه دانشی برای مفاهیم استفاده می‌کند. هسته اولیه دارای مفاهیم و عملگرهای بنیادی است که برای اکتساب خودکار دانش مورد نیازند. به عبارت دیگر هسته حاوی فرادانش لازم جهت ساخت انواع هستان شناسی‌ها بر روی آن، با توجه به دانش دریافتی از محیط از طریق جملات ساده زبان فارسی است. سیستم با شروع بکار با این هسته، با تراکنش‌هایی که به زبان طبیعی با محیط خود دارد، لایه‌های بعدی هستان شناسی را حول هسته اولیه می‌سازد. تنها وسیله ارتباطی هستی با محیط، متون نوشتاری زبان فارسی و هدف آن ساخت هستان شناسی براساس اطلاعات دریافتی از متون ورودی است. تأمین این اطلاعات وظیفه «مأمور پردازش زبان طبیعی» است که واژگان یکی از اجزاء اصلی در ساختار پایگاه دانش آن است. شکل (۱) ساختار کلی سیستم هستی را نشان می‌دهد.

#### ۲-۶- ویژگی‌های واژگان در هستی

ساختار واژگان برای حفظ دانش لغات در سیستم طراحی شده است. این ساختار در ابتدا حاوی بخشی از دانش اولیه در مورد لغات است و سپس در طول حیات سیستم<sup>۱۲</sup> دانش‌های جدید توسط سیستم، اخذ [۹] و در آن درج

شکل (۱) ساختار کلی سیستم.





شکست: اسم). با توجه به اینکه در زبان منتخب (زبان لیسپ [۱۱]) امکان تخصیص پویای حافظه در صورت لزوم وجود دارد و اطلاعات به صورت نمادی پردازش می‌شوند، لذا می‌توانیم برای مدخل‌های مختلف تعداد و نوع ویژگی‌ها را متغیر در نظر بگیریم و هر گونه اطلاعات مرتبط با یک واژه را تنها در یک مدخل در کنار هم نمایش دهیم.

- نحوه نگاشت مدخل‌ها به واحدهای لغوی: با توجه به نحوه برخورد با کلمات هم‌شکل، نحوه نگاشت مدخل‌ها به واحدهای لغوی به صورت یک به یک خواهد بود. چرا که هر مدخل به یک واحد لغوی (یک کلمه با نمود یا مفهوم‌های متفاوت) تخصیص می‌یابد.

- نحوه تعیین واحدهای لغوی: واحدهای لغوی در این واژگان کلمات هستند. زیر کلمات (پیشوندها، پسوندها، میانوندها) در ساختار جداگانه‌ای در بخش ساختار واژه نگه‌داری می‌شوند. زیر کلمات (کلمات مرکب یا چند کلمه‌ای‌ها) با توجه به انتخاب کلمات ساده در متون آزمون، در واژگان قرار نمی‌گیرند.

- یگانگی ساختار واژگان: واژگان مورد نظر از دو ساختار جداگانه بهره می‌برد: ساختار واژه‌ها و ساختار واژه‌سازها. ساختار واژه‌ها (که در این مقاله واژگان نامیده می‌شود) حاوی اطلاعات لازم در مورد واژه‌های زبان است و دارای ساختار یگانه است. به عبارت دیگر ساختار آن برای گروه‌ها و انواع مختلف واحدهای لغوی (اعم از اسم، فعل و ...) یکسان است. اما ساختار واژه‌سازها حاوی اطلاعات ساختار واژه‌ها در مورد وندهاست و اگرچه خود ساختار یگانه دارد ولی با ساختار واژگان متفاوت است. این ساختار از این پس واژگان ساختار واژه نامیده می‌شود.

- تعیین مدخل‌های اصلی و زیر مدخل‌ها: در این واژگان مدخل‌ها تک سطحی بوده و زیر مدخلی وجود ندارد.

- ترتیب واحدهای با تشابهات صوری: واحدهای با تشابهات صوری به عکس ترتیب ورود (یا یادگیری) حفظ می‌شوند.

- تکنیک انجام توضیحات: از آنجا که کاربر، سیستم درک و یادگیری است، توضیحات به صورت متن نوشتاری تنها برای مواردی که لازم باشد عیناً در خروجی برای کاربر انسانی سیستم چاپ شود، ذخیره می‌شود. خود سیستم از این بخش استفاده‌ای نمی‌کند.

ب- اهم اطلاعات ذخیره شده برای هر مدخل در واژگان در هستی با توجه به اینکه در حال حاضر واژگان جهت درک متون با کلمات و جملات ساده زبان بکار می‌رود، لذا از ذخیره سازی اطلاعات واجی و همچنین اطلاعاتی که جهت تولید زبان بکار می‌روند (مانند برخی اطلاعات جانشینی) و اطلاعات لازم جهت ساخت کلمات مرکب (برخی اطلاعات اشتقاقی)

صرف نظر می‌کنیم. با توجه به این موارد اطلاعات ذخیره شده در واژگان هستی به شرح زیر خلاصه می‌شود:

- شکل نوشتاری کلمه

- دانش ساختار واژه: در هستی بخش اعظم دانش ساختار واژه زبان در قالب قواعد ساختار واژه‌ها و یا به صورت جداگانه در ساختار واژگان ساختار واژه (حاوی واژه‌سازها مانند پیشوندها، میانوندها و پسوندها) ذخیره می‌شود. بخش دیگری از این دانش که خاص کلمات و واژه‌هاست، در واژگان در مقابل هر واژه آورده می‌شود. مانند نحوه جمع بستن، جمع مکسر یا شکل مفرد یک اسم، ریشه کلمه، شکل ماضی یا مضارع فعل و ... به عبارت دیگر ما بیشتر به حفظ اطلاعات تصرفی توجه داریم و فعلاً اطلاعات اشتقاقی را نگهداری نمی‌کنیم.

- دانش نحوی: این بخش شامل طبقه یا نوع کلمه و ویژگی‌ها یا زیر طبقه کلمه می‌باشد. این زیر طبقات ویژگی‌های جزئی تری از کلمه را در طبقه خود بیان می‌کنند. به عنوان نمونه در طبقه فعل زیر طبقات می‌توانند نوع فعل (مثلاً ماضی یا مضارع) و یکسری ویژگی‌های دودویی دیگر مانند متعدی یا مجهول بودن و در طبقه اسم شامل ویژگی‌های دودویی مانند جمع یا مفرد، ذات یا معنی، شمارا یا ناشمارا، عام یا خاص باشند. البته ما در مرحله اول پروژه با اطلاعات هم‌نشینی کار نمی‌کنیم ولی در مراحل بعد برخی اطلاعات هم‌نشینی کلمه مانند نحوه ترکیب آن با کلمات دیگر و یا با حروف اضافه و طبقات نحوی مجاز همراه شونده با فعل در این قسمت درج خواهد شد.

- دانش معنایی: در هستی، معنای اصلی واژه توسط اشاره به مفهوم متناظر آن در هستان شناسی تعیین می‌شود ولی ویژگی‌های خاص معنایی و یا محدودیت‌های گزینشی خاص واژه در واژگان قید می‌گردند. لذا بسیاری از اطلاعات معنایی که بعضاً در واژگان‌های محاسباتی دیده می‌شوند ([۵]) مانند اطلاعات جانشینی، در هستان شناسی هستی تعبیه می‌شود و واژگان از طریق بخش معنایی واژه به این اطلاعات دسترسی دارد.

- دانش کاربردی: در هستی بیشتر دانش کاربردی مورد نظر (مانند بخشی از اشارات ضمنی، پیش شرایط و نتایج و آثار اعمال و ...) در هستان شناسی جای دارد و بخش دیگری که متصل به واژه است (مانند ارزش واژه در میان واژه‌های هم‌گروهش) در واژگان و در ارتباط با هستان شناسی ذخیره می‌گردد. مثلاً برای ارزش گذاری واژه‌های داغ، گرم و ولرم در واژگان ارزش نسبی حرارتی به این واژه‌ها داده می‌شود که این ارزش مرتبط با مفهوم دما یا حرارت در هستان شناسی است و لذا بین همه واژه‌های مرتبط با این مفهوم مشترک بوده و قابل قیاس است.



- ویژگی های واژه مرکب: طبقه نحوی و ویژگی های معنایی واژه حاصل از ترکیب را بیان می کند.

به عنوان نمونه، واژه ساز «بان» با واژه هایی با طبقه نحوی «اسم» می تواند ترکیب پسوندی بسازد و حاصل، واژه ای با طبقه نحوی «اسم» و ویژگی معنایی «محافظة» خواهد بود.

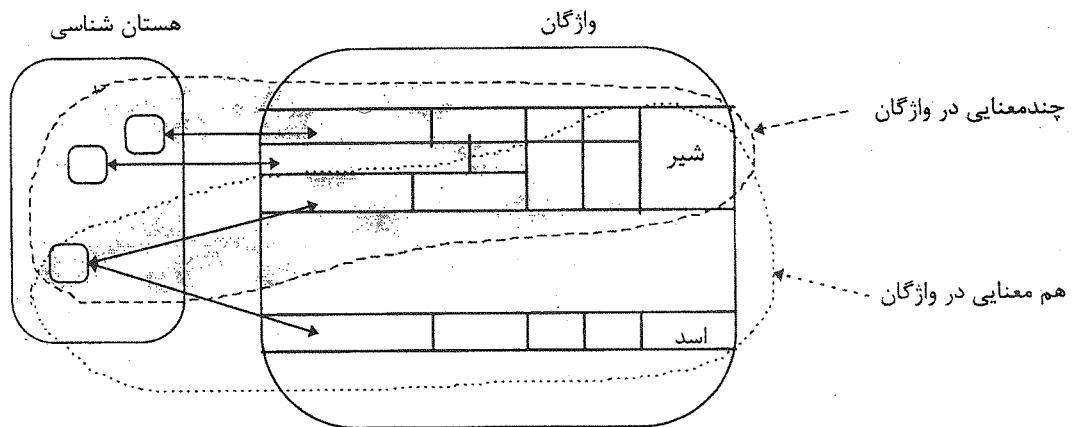
### ۵-۶- ارتباط واژگان با سایر اجزاء پایگاه دانش

در واژگان هستی، بخش اعظم اطلاعات معنایی و کاربردی از طریق ارتباط با هستان شناسی تأمین می گردد. برای هر کلمه ممکن است بیش از یک ارتباط با هستان شناسی برقرار باشد (چند معنا بودن)<sup>17</sup> مثل کلمه «شیر» که بیش از یک معنا دارد و لذا با بیش از یک ارتباط معنایی به هستان شناسی مرتبط می شود. همچنین ممکن است به یک مفهوم در هستان شناسی بیش از یک کلمه مرتبط باشد (هم معنا بودن)<sup>18</sup>. مثل مفهوم "speak" که معنایی برای کلمات مترادف

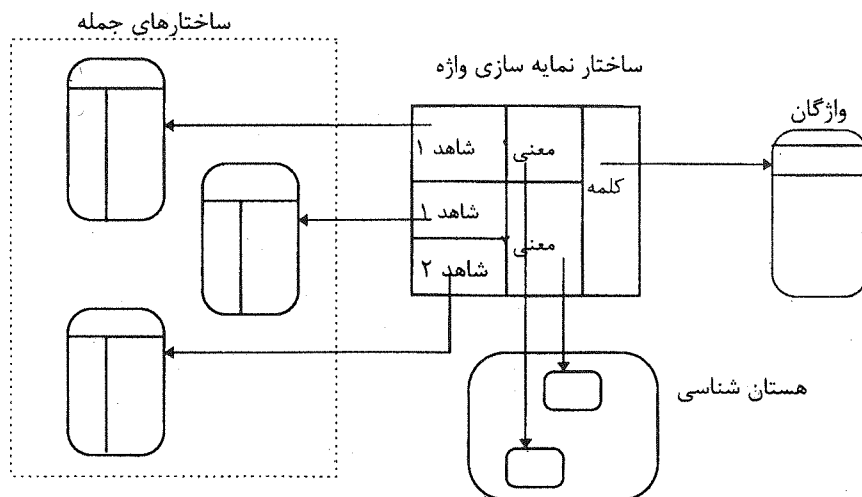
گفتن، بیان کردن و صحبت کردن است. شکل (۳) طرح واره ای از ارتباط واژگان و هستان شناسی را در هستی نشان می دهد.

از آنجا که سیستم معانی لغات (مفاهیم و ارتباطات) را بتدریج و از روی متون دریافتی یاد می گیرد، به دلایل مختلف ممکن است لازم باشد در آموخته های خود تجدید نظر کرده و یا آنها را تغییر و تصحیح نماید. به این منظور لازم است سابقه ای از مسیر یادگیری خود را حفظ کند. ساختار نمایه سازی واژه ها، ساختاری است که با همین هدف طراحی شده است. این ساختار میان واژگان، هستان شناسی و ساختار متن قرار می گیرد و همه مفهوم های کلمه را در متن مشخص می کند. این ساختار برای

هر واژه لیستی از شاخص های ساختار جملات متن نگه می دارد که در آن جمله، آن کلمه با معنی مورد نظر بکار رفته است. این ساختار در حافظه کاری<sup>۱۹</sup> سیستم نگهداری و با



شکل (۳) طرح واره ای از ارتباط واژگان و هستان شناسی.



شکل (۴) طرح واره ای از ارتباط واژگان، هستان شناسی و ساختارهای جمله در یک ساختار متن.

تعویض متن ورودی، آخرین نتایج مندرج در آن به ساختارهای واژگان و هستان شناسی منتقل و ساختار خالی خواهد شد. شکل (۴) طرح واره ای از ارتباط واژگان، هستان شناسی و ساختارهای جمله را در یک ساختار متن نشان می دهد.

## ۶-۶- یادگیری لغات و معانی تازه

واژگان هستی در ابتدای کار دارای تعداد محدودی واژه است که معانی بیشتر آنها را نمی داند. سیستم به مرور با دریافت جملات ساده، با استفاده از قواعد دستور زبان فارسی و با توجه به جملات دیگر متن، اطلاعات مربوط به کلمات موجود را تکمیل و یا کلمات جدید را به واژگان می افزاید. برای روشن تر شدن بحث، عمل یادگیری سیستم را برای مواجهه با کلمه جدید، و تکمیل اطلاعات ساختواژی، نحوی و معنایی به تفکیک بررسی می کنیم.

**الف - کلمه جدید:** در حالتی که یکی از کلمات متن در واژگان موجود نباشد ابتدا امکان تبدیل آن به یکی از کلمات موجود در واژگان با توجه به قواعد ساختواژی موجود و بر اساس واژگان ساختواژی بررسی می شود. مثلاً با برخورد به کلمه «کتاها» در شرایطی که واژه «کتاب» در واژگان هست، گذر از یک مرحله تحلیل ساختواژی ما را با حذف علامت جمع به یک کلمه شناخته شده می رساند. اما اگر نتوانستیم کلمه را ساده تر کنیم (یعنی یا اجزاء کلمه در واژگان نبود یا کلمه قابل تفکیک نبود)، با یک کلمه جدید مواجهیم. در این حالت مدخلی برای آن در واژگان در نظر گرفته شده، بدنبال تکمیل اطلاعات این مدخل خواهیم بود.

**ب) اطلاعات ساختواژی:** در مرحله اول ساخت هستی ما از میان اطلاعات ساختواژی به نحوه جمع بستن، جمع مکسر یا شکل مفرد یک اسم جمع، ریشه کلمه، و شکل ماضی یا مضارع فعل توجه داریم. جمع های مکسر و اسامی جمع در بسیاری موارد از طریق مطابقت فاعل و فعل قابل تشخیص می باشند. همچنین نحوه جمع بستن یک اسم و یا ریشه کلمه با گذر واژه از تحلیل گر ساختواژی استخراج می شود. اما در مورد بقیه موارد فوق در مرحله اول به یادگیری از طریق محاوره با کاربر بسنده کرده ایم.

**ج - اطلاعات نحوی:** برای دریافت و تکمیل اطلاعات نحوی یک واژه دو مسیر در سیستم وجود دارد. اول اینکه هنگام عبور واژه از مرحله تحلیل ساختواژی ممکن است نقش نحوی واژه قابل تشخیص باشد چرا که حاصل اعمال بسیاری قواعد ساختواژی برای سیستم شناخته شده است. مثلاً اگر سیستم به کلمه ناشناس باغبان برخورد کند، در شرایطی که پسوند «بان» در ساختار واژه سازها آمده باشد، پس از عبور از تحلیلگر ساختواژی نقش نحوی کلمه، اسم تشخیص داده

می شود، چرا که سیستم می داند پسوند «بان» بر سر اسم آمده و اسم دیگر می سازد.

مسیر دوم دریافت اطلاعات نحوی، از طریق تجزیه نحوی جمله است. در این حالت ابتدا سیستم با کمک یک تجزیه گر چارت بالا به پائین<sup>20</sup> به دنبال قاعده ای در قواعد دستور زبان می گردد که با جمله جاری تطابق داشته باشد. این قاعده نقش دستوری کلمه یا کلمات ناشناس را تعیین می کند. در صورتی که بیش از یک قاعده برای جمله جاری یافت شده و یا بیش از یک نقش برای کلمه ناشناس بدست آمد، همه حالات ممکن برای تأیید یا رد بعدی در ساختار فرضیات حفظ خواهند شد. تأیید یا رد نتایج حاصله بر اساس جملات بعد و یا تراکنش با کاربر و توسط مأمور یادگیری انجام می پذیرد (برای بحث مفصل تر در مورد شرح وظایف و نحوه عمل مأمور یادگیری رجوع شود به [۱۳]). همچنین دریافت اطلاعاتی در مورد زیر طبقات افعال نیز در همین مرحله با توجه به کلمات یک جمله میسر است که در مراحل بعدی پروژه پیاده سازی خواهد شد.

**د - اطلاعات معنایی و کاربردی:** هنگامی که جمله با یکی از قواعد دستوری زبان تطابق یافت، با استفاده از الگوهای معنایی در نظر گرفته شده برای قواعد دستوری، نقش<sup>21</sup> کلمات جمله استخراج شده و در ساختاری به نام ساختار جمله حفظ می شود. ساختار جمله در حقیقت قالب حالتی<sup>22</sup> است که برای جملات کنشی شامل نقش های کنش، کنش گر، کنش پذیر، زمان، مکان، ابزار، توصیف گر کنش، مبدأ و مقصد است و برای جملات حالتی حاوی حالت پذیر، حالات، زمان و مکان می باشد.

پس از ساخته شدن ساختار جمله، به ازاء هر کلمه ای در جمله که معنای آن ناشناخته است مفهومی<sup>23</sup> در هستان شناسی ایجاد می گردد. این مفاهیم با توجه به طبقه کلمه و نقش آن در جمله در بالاترین سطح زیر هسته به ساختار هستان شناسی متصل می شوند<sup>24</sup> (به عنوان نمونه، اسامی تحت مفهوم Object، افعال تحت مفهوم Action و صفات تحت مفهوم Property). سپس این مفهوم از محل پدر اولیه خود به طرف برگ های سلسله مراتب ارث بری در هستان شناسی حرکت می نماید تا خاص ترین پدر خود را بیابد. منظور از یافتن خاصترین پدر، یافتن مفهومی است که اولاً از نظر ویژگی با مفهوم جدید تطابق داشته باشد، ثانیاً در میان مفاهیم مطابق، در پائین ترین سطح ممکن در سلسله مراتب قرار گرفته باشد. عمل یافتن تطابق خود رویه پیچیده ای دارد. ابتدایی ترین شکل تطابق در این مرحله، عدم وجود تضاد است.<sup>25</sup> در صورتی که همه یا هیچیک از فرزندان یک پدر با مفهوم جدید مطابقت داشته باشند، عمل انتقال متوقف و

مفهوم تحت پدر قرار می گیرد. در صورتی که بیش از یک فرزند مطابقت داشته باشد، با استفاده از یک تکنیک نگاه به جلو<sup>26</sup>، تا N سطح جلوتر (در مرحله اول N=1) آزمایش می شود تا نزدیک ترین تطابق پیدا شود. در صورتی که همچنان بیش از یک انتخاب وجود داشت، در این مرحله مفهوم به گره پدر منتسب می شود<sup>27</sup>.

به مرور که مفاهیم مختلف به ازاء کلمات زبان تشکیل می گردند، این مفاهیم بر اساس ویژگی های مشترک در هم ادغام و یا بر اساس عوامل تفکیک به دسته های ریزتر شکسته می شوند. در شرایطی که اختلافی میان معنی موجود و معنی دریافتی از جمله جاری بیاید، نخست معنی جدید را با جملات قبل مقایسه می کند. در صورتی که معنی جدید با جملات قبل تطابق داشته باشد (مثلاً زیر مجموعه یا زبر مجموعه معانی سابق باشد) معنی واژه را به معنی جدید تصحیح می کند. در غیر این صورت معنی جدید را به عنوان مفهوم دیگری از واژه مورد نظر می آزماید. در هر دو صورت پیغامی مبنی بر تصمیم سیستم به کاربر انسانی ارسال می شود تا کاربر در صورت تمایل این تصمیمات را تصحیح یا تأیید نماید.

مثلاً فرض کنید اولین جمله ورودی جمله «حسن سیب را خورد» بوده، واژگان حاوی مدخل هایی با محتویات زیر باشد:

حسن: اسم خاص، مفرد.

سیب: اسم، مفرد،

را: حرف اضافه - نشانه

خورد: فعل ماضی، متعدی

ساختار جمله ساخته شده بر اساس قاعده دستوری:

«جمله ← گروه اسمی (کنش گر) + گروه اسمی (کنش پذیر) + (را + فعل متعدی (کنش))»

نشان می دهد که حسن فاعل (کنش گر)، سیب مفعول (کنش پذیر) و خورد کنش جمله است. لذا در هستان شناسی مفاهیمی برای حسن، سیب و خورد تحت ابر مفاهیم چیز و کنش ایجاد می شود. در این مرحله چون این ابر مفاهیم خود در برگ ها قرار دارند، انتقالی صورت نمی گیرد. سپس روابط قوه کنش کاری، قوه کنش پذیری، قوه کنش کاری بر وقوه کنش پذیری از، میان این مفاهیم برقرار می گردد. پس از آن با رسیدن به جمله ای مثل «حسین سیب را خورد» با توجه به وجود مفاهیمی برای سیب و خورد، مفهومی برای حسین در هستان شناسی تشکیل می شود. در این لحظه حسن و حسین هر دو توان کنش گری یک کنش بر روی یک کنش پذیر را دارند و این یک اشتراک محسوب می شود. به مرور که جملات بیشتری توسط سیستم دریافت می شود، این اشتراکات افزایش می یابد. هنگامی که اشتراکات دو مفهوم از یک

آستانه ای (آستانه ادغام<sup>28</sup>) فراتر رفت و در ضمن میزبان اختلافاتشان از آستانه دیگری (آستانه تفکیک<sup>29</sup>) کمتر بوده، تضاد ویژگی نیز میان نشان وجود نداشت، می توان این دو مفهوم را به یک مفهوم کلی تر ادغام و یا ابر مفهوم<sup>30</sup> مشترکی برای آنها در نظر گرفت ([۱۲]).

## ۶-۲. پیاده سازی

واژگان در این سیستم به زبان لیست و در محیط Allegro CL 6.0 تحت سیستم عامل Windows پیاده شده است. جهت ذخیره و بازیابی اطلاعات در واژگان از ساختار Hash-table استفاده شده و جستجو در آن مرحله اول با توابع hash موجود در محیط انجام می گیرد. واژگان در مرحله اول حاوی حروف اضافه، افعال ربطی، و تعدادی (کمتر از ۱۰ مورد) اسم و فعل است که به صورت تصادفی انتخاب شده اند.

## ۷. نتیجه گیری

در گذشته سیستم های پردازش زبان طبیعی (NLP)، به واژگان به عنوان یک جزء لازم که از توجهات نظری اندکی برخوردار بود و به سادگی با کمک مجموعه ای از برنامه ها ساخته می شد، نگاه می کردند. در حالیکه امروزه مرکزیت واژگان در سیستم های پردازش زبان مورد تأکید است. انتخاب صحیح ویژگی های ساختاری و محتوایی و تعیین دقیق نوع و ساختار واژگان مورد نیاز بر اساس کاربرد سیستم مورد نظر از وظایف طراح است که نقش اساسی در عملکرد سیستم خواهد داشت. در نمونه عملی مورد مطالعه اگر چه هدف اصلی سیستم، پردازش زبان نبوده و عمل پردازش زبان به صورت محدود (بر روی ورودی محدود) و تنها به عنوان تغذیه کننده سیستم اصلی انجام می شود ولی دلیل تعیین کننده بودن واژگان در کارایی سیستم نهایی و اکتساب دانش مورد نیاز از متون به ساختار واژگان توجه خاص نموده ایم. این ساختار به سهولت قابل تعمیم برای دربرگیری ورودی های نامحدود زبان و یا تعمیم به کاربرد تولید متن می باشد.

## زیر نویس ها

1- Machine Readable Dictionaries

2- Corpora

۳ - منابع زبانی به دو دسته اصلی تقسیم می شوند: الف - منابع اولیه منابع زبانی مورد استفاده هستند که به صورت عام وجود دارند و خاص استخراج اطلاعات واژگانی طراحی نشده اند. از جمله این منابع می توان به متون نوشتاری یا بخش های ضبط شده ای از گفتار اشاره

نمود. منابع ثانویه: منابعی هستند که خود از منابع اولیه یا ثانویه دیگر استخراج شده‌اند مانند دیکشنری‌ها یا دستور زبان‌های موجود.

- 4- Dictionary
- 5- Lexical Unit
- ۶- Morphemes - اجزاء ساخت یک کلمه، مثلاً «نا» (نقی) + «مادر» + «ی» تکرارهای ساخت کلمه «نامادری» هستند.
- ۷- Sememes - اجزاء معنا که در کاربردهای مختلف متفاوت است مثلاً مفاهیم اولیه متناظر در هستان شناسی
- 8- Entry
- 9- Paradigmatic
- 10- Syntagmatic
- 11- Analogical
- 12- Diasystematic
- 13- Ontology
- ۱۴ - از لحظه آغاز بکار سیستم تا هر زمانی که سیستم مورد استفاده قرار گیرد.
- 15 - Sense
- 16- Sense Group
- 17- Polysemy
- 18- Synonymy
- ۱۹ - حافظه کاری محل نگهداری اطلاعات میان اکتسابی از محیط و حفظ ساختارهای کمکی سیستم برای انجام تحلیلات لازم است.
- 20- Top Down Chart Parser
- 21- Role
- 22- Case Frame
- 23- Concept

## مراجع

- ۲۴ - هسته هستان شناسی حاوی مجموعه کوچکی از مفاهیم و عملگرهای ابتدایی برای شروع به کار سیستم است. این مفاهیم در مرحله اول شامل: Thing در بالاترین سطح و Object, Action, Property در سطح دوم می‌باشند. همچنین عملگرهای اولیه در این مرحله عملگرهای Add جهت درج یک مفهوم، Delete جهت حذف یک مفهوم، Update جهت به‌نگام سازی ویژگی‌های یک مفهوم، Merge جهت ادغام دو یا چند مفهوم، Split جهت تفکیک یک مفهوم و Transfer جهت انتقال یک مفهوم در ساختار هستان شناسی، می‌باشند.
- ۲۵ - به منظور یافتن تطابق، هر ویژگی دارای نشانگری است که نشان می‌دهد مقدار آن قابل ادغام، اجتماع، بازنویسی و یا تغییر توسط مقدار جدید هست یا نه. به منظور یافتن تطابق لازم است ویژگی‌های غیر همسان، اولاً متضاد نباشند، ثانیاً مقادیر آنها متضاد نباشند، ثالثاً این مقادیر براساس نشانگرشان با هم قابل تطابق باشند. تشخیص تضاد براساس ویژگی‌های منصوب به مفهوم ویژگی مورد نظر در هستان شناسی میسر می‌گردد.
- 26- Look Forward
- ۲۷ - توضیح در مورد محتوای اولیه هستان شناسی، رویه‌های ادغام، تفکیک و انتقال مفاهیم و نحوه تشخیص تضاد، همسانی و تطابق، بحث مفصلی را می‌طلبد که باتوجه به تأکید این مقاله بر بخش واژگان، در حوصله این بخش نیست. شرح این موارد در مرجع [۱۲] آورده شده است.
- 28- Marge Treshold
- 29- Split Treshold
- 30- Super Concept

- [1] Atwell E. S., "Machine Readable Dictionaries and Lexical Databases", School of computer Studies, Leeds university, 1998.
- [2] Malmkjaer K., "The Linguistics Encyclopedia", Routledge, 1996.
- [۳] شمس فرد م.، «واژگان: مباحث زبانی و محاسباتی»، گزارش فنی شماره ۱۰۱-۱۰۲-۴۰۰۱-۷۸، آزمایشگاه سیستم‌های هوشمند، دانشکده مهندسی کامپیوتر، دانشگاه صنعتی امیرکبیر، شهریور ماه ۱۳۷۸.
- [4] Cavazza M., Zweigenbaum P., "Lexical Semantics: Dictionary or Encyclopedia?", in [14], pp 336-347, 1995.
- [5] Miller G. A., "WordNet: A Lexical Database for English", Communication of the ACM, Vol 38, No. 11, p 39-41, NOV. 1995.
- [6] Nirenburg S., et al, "Lexicons in the Mikrokosmos project", AISB Workshop on Multilinguality in the Lexicon, Brighton, UK, 1996.
- [۷] عبدالله زاده ا.، شمس فرد م.، «هستی مدل پیشنهادی هستان شناسی در سیستم‌های درک زبان طبیعی»، هفتمین کنفرانس مهندسی برق ایران، اردیبهشت ۱۳۷۸
- [8] Shamsfard M., Abdollahzadeh A., "A Basis for Evolutinary Ontology Construction", 18th Iasted International Conference on Applied Informatics

- (AI 2000), Feb. 2000, Austria.
- [۹] شمس فرد م.، «یادگیری ماشینی»، گزارش فنی شماره ۱۰۲-۱۰۱-۴۰۰۱-۷۸، آزمایشگاه سیستم‌های هوشمند، دانشکده مهندسی کامپیوتر، دانشگاه صنعتی امیرکبیر، شهریور ماه ۱۳۷۸.
- [۱۰] شمس فرد م.، «طراحی محیط پیاده سازی و تست هستی: سیستم یادگیری هستان شناسی مبتنی بر درک متون فارسی»، گزارش فنی شماره ۱۰۴-۱۰۱-۴۰۰۱-۷۹، آزمایشگاه سیستم‌های هوشمند، دانشکده مهندسی کامپیوتر، دانشگاه صنعتی امیرکبیر، خرداد ماه ۱۳۷۹.
- [11] Winston P. H., Horn B. K. P. "Lisp", 3rd edition, Addison Wesley, 1993.
- [۱۲] شمس فرد م.، «هسته هستان شناسی، مفاهیم و عملگرهای اولیه»، گزارش فنی شماره ۱۰۶-۱۰۱-۴۰۰۱-۸۰، آزمایشگاه سیستم‌های هوشمند، دانشکده مهندسی کامپیوتر، دانشگاه صنعتی امیرکبیر، بهار ۱۳۸۰.
- [۱۳] شمس فرد م.، «طراحی ساختار اولیه هستی: سیستم یادگیری هستان شناسی مبتنی بر درک متون فارسی»، گزارش فنی شماره ۱۰۵-۱۰۱-۴۰۰۱-۷۹، آزمایشگاه سیستم‌های هوشمند، دانشکده مهندسی کامپیوتر، دانشگاه صنعتی امیرکبیر، شهریور ماه ۱۳۷۹.
- [14] Saint-Diezer P., Viegas E., "Computational Lexical semantics", Cambridge university press, 1995.