

# پژوهشگاه علوم انسانی و مطالعات فرهنگی

## پژوهشکده زبانشناسی

پایان نامه

برای دریافت درجه کارشناسی ارشد زبانشناسی همگانی

موضوع

برچسب‌دهی دستوری پیکره زبان فارسی

پژوهشی در زبانشناسی رایانه‌ای

استاد راهنما: دکتر سید مصطفی عاصی

استاد مشاور: دکتر شاپور اعتماد

پژوهش از: محمد حاجی عبدالحسینی

اسفند ۱۳۷۵

## فهرست مندرجات

| یک | پیشگفتار  |
|----|---|
| ۱  | فصل ۱ زیانشناسی پیکره بنیاد                                       |
| ۱  | ۱-۱ درآمد   |
| ۲  | ۲-۱ نشانه گذاری پیکره   |
| ۳  | ۱-۲-۱ نشانه گذاری نگارشی  |
| ۴  | ۲-۲-۱ نشانه گذاری آوایی/واجی                                      |
| ۴  | ۳-۲-۱ نشانه گذاری زیر زنجیری                                      |
| ۵  | ۴-۲-۱ برجسب دهی دستوری  |
| ۵  | ۵-۲-۱ تقطیع نحوی  |
| ۵  | ۶-۲-۱ نشانه گذاری معنایی  |
| ۶  | ۷-۲-۱ نشانه گذاری کارکردشناختی/کلامی                              |
| ۶  | ۳-۱ چکیده   |
| ۷  | فصل ۲ برجسب دهی دستوری  |
| ۷  | ۱-۲ تاریخچه نشانه گذاری نحوی                                      |
| ۹  | ۲-۲ روش   |
| ۱۱ | ۳-۲ چکیده   |
| ۱۲ | فصل ۳ برجسب دهی دستوری توزیعی                                     |
| ۱۳ | ۱-۳ استقرا تنها براساس گونه های واژه ها                           |
| ۱۴ | ۲-۳ استقرا براساس گونه های واژه ها و بافت                         |
| ۱۵ | ۳-۳ استقرا براساس گونه های واژه ها و بافت منحصر به بافت های طبیعی |
| ۱۵ | ۴-۳ استقرا براساس بردارهای تعمیم یافته                            |
| ۱۵ | ۵-۳ نتایج آزمایش های شوتس   |

|    |  |
|----|--|
| ۱۶ | ۶-۳ چکیده                                      |
| ۱۷ | فصل ۴ پروژه برچسب‌دهی دستوری پیکره فارسی       |
| ۱۷ | ۱-۴ معرفی پروژه                                |
| ۱۸ | ۲-۴ مجموعه برچسب                               |
| ۱۹ | ۱-۲-۴ معیارهای ایجاد یک طرح نشانه‌گذاری        |
| ۲۰ | ۱-۱-۲-۴ از دید نشانه‌گذار                      |
| ۲۰ | ۲-۱-۲-۴ از دید کاربر                           |
| ۲۰ | ۳-۱-۲-۴ معیارهای زیانشناختی خارجی              |
| ۲۲ | ۲-۲-۴ برچسب‌های به کار رفته در پروژه حاضر      |
| ۲۵ | ۳-۴ برنامه‌های موجود در بسته نرم‌افزار کاربردی |
| ۲۶ | ۱-۳-۴ VGAF.COM                                 |
| ۲۷ | ۲-۳-۴ PROISAMD.EXE                             |
| ۲۸ | ۳-۳-۴ TAG.BAT                                  |
| ۲۸ | ۴-۳-۴ TAGGER.EXE                               |
| ۳۲ | ۱-۴-۳-۴ روال اصلی                              |
| ۳۳ | ۲-۴-۳-۴ زیر- روال ContextFinder                |
| ۳۵ | ۳-۴-۳-۴ زیر- روال FindCon                      |
| ۳۸ | ۴-۴-۳-۴ زیر- روال FindGeneralCon               |
| ۳۸ | ۵-۴-۳-۴ زیر- روال ViewCon                      |
| ۳۹ | ۶-۴-۳-۴ زیر- روال MDBMaker                     |
| ۴۰ | ۷-۴-۳-۴ زیر- روال Categorizer                  |
| ۴۵ | ۸-۴-۳-۴ زیر- روال MDBCreate                    |
| ۴۵ | ۹-۴-۳-۴ زیر- روال WordTypeMDB                  |
| ۴۷ | ۱۰-۴-۳-۴ زیر- روال GeneralConMDB               |

|    |   |
|----|---|
| ۴۷ | AskTag ۱۱-۴-۳-۴ زیر- روال                                     |
| ۴۸ | ManualTag ۱۲-۴-۳-۴ زیر- روال                                  |
| ۴۹ | TagTXT ۱۳-۴-۳-۴ زیر- روال                                     |
| ۵۲ | FastTXT ۱۴-۴-۳-۴ زیر- روال                                    |
| ۵۳ | TXTSET.EXE ۵-۳-۴  |
| ۵۴ | YCHANGE.EXE ۶-۳-۴   |
| ۵۴ | NDXSETUP.EXE ۷-۳-۴  |
| ۵۵ | ۴-۴ چکیده   |
| ۵۷ | فصل ۵ آزمایش بسته نرم افزار و پیشنهادهایی برای پروژه های بعدی |
| ۵۷ | ۱-۵ برچسب دهی متن اول   |
| ۵۷ | ۲-۵ برچسب دهی متن دوم   |
| ۵۸ | ۱-۲-۵ برچسب دهی دستی هزار واژه اول فایل نمایه                 |
| ۵۸ | ۲-۲-۵ طبقه بندی تنها با استفاده از میزان شباهت بردارها        |
| ۵۹ | MSM= ۱۰۰ ۱-۲-۲-۵  |
| ۶۰ | MSM= ۸۰ ۲-۲-۲-۵   |
| ۶۲ | MSM= ۶۰ ۳-۲-۲-۵   |
| ۶۲ | ۳-۲-۵ طبقه بندی با استفاده از بردارهای حاوی برچسب واژه ها     |
| ۶۳ | ۳-۵ نتیجه گیری  |
| ۶۳ | ۴-۵ پیشنهادهایی برای پروژه های بعدی                           |
| ۶۴ | ۵-۵ چکیده   |
| ۶۵ | کتابنامه  |

## پیشگفتار

گسترش سریع و فزاینده رایانه و بسط بی حد و حصر امکانات این ابزار قدرتمند، موجب گردیده است که امروزه تقریباً تمام شاخه‌های علوم و هنر از آن بهره‌مند گردند. زبانشناسی نیز از این دسته مستثنی نیست. اتفاقاً زبان جزء اولین حوزه‌هایی بود که کاربرد رایانه در آن مد نظر قرار گرفت. منظور دوره تب ساختن ماشین‌های ترجمه در دهه‌های ۱۹۵۰ و ۱۹۶۰ است. صحیح است که پیچیدگی زبان و سادگی رایانه‌های ابتدایی و نیز ضعف نظریه‌های زبانی آن دوران موجب رکود فعالیت‌های زبانشناسی رایانه‌ای گردید، اما از دهه ۱۹۷۰ این گونه فعالیت‌ها و مطالعات با بینشی نواز سرگرفته شد. در این زمان زبانشناسی رایانه‌ای به حدی گسترش یافته است که صحبت از مترجم‌های شفاهی ماشینی، سیستم‌های یادگیرنده، تولید متن به وسیله رایانه‌های بسیاری فعالیت‌های دیگر به میان آمده است.

اما زبانشناسی رایانه‌ای در کشور ما هنوز بسیار ناشناخته باقی مانده و چشم‌اندازهای وسیع نظری و کاربردی آن از دید عده کثیری به دور است. نگارنده به دلیل علاقه فراوان به شاخه زبانشناسی رایانه‌ای بر آن شد تا پایان‌نامه کارشناسی ارشد خود را در این زمینه انتخاب نماید. این کار هم می‌توانست به آشنایی بیشتر نگارنده به این رشته و هم به شناساندن آن به جامعه زبانشناسی ایران کمکی هر چند ناچیز کرده باشد. استاد راهنمای محترم نگارنده، جناب آقای دکتر عاصی، ضمن ارائه رهنمودهای گوناگون، پایگاه داده‌های زبان فارسی و اهداف آن را معرفی نموده و پیشنهاد کردند که چه بهتر موضوع یک پایان‌نامه تنها پایان کار نبوده بلکه کاربردهای عملی در یک پروژه وسیع‌تر داشته باشد و چه قدر این امر صحیح است. بدین ترتیب، نگارنده موضوع برجسب‌دهی دستوری پیکره زبان فارسی را که یکی از مراحل ایجاد و گسترش دامنه کاربرد یک پیکره زبانی است انتخاب نمود.

هدف از یک طرح برجسب‌دهی، مشخص کردن مقوله دستوری هر یک از اقلام واژگانی موجود در یک متن است. در گذشته چنین کاری به طور دستی انجام می‌گرفت، اما اکنون امکان انجام خودکار آن با استفاده از رایانه پیدا شده است. البته طبق معمول

پیچیدگی‌های زبان موجب می‌گردد که کار آن طور که به نظر می‌رسد ساده نباشد. به همین دلیل است که، روش‌های متفاوتی برای انجام این کار در کشورهای دیگر ابداع شده است.

از آنجا که این کار برای اولین بار روی زبان فارسی انجام می‌پذیرد و برای آن که بتوان سیستم را در مدت زمان تعیین شده برای انجام پایان‌نامه فعال نمود، تصمیم گرفتم طرحی ساده و نیمه خودکار را برای این کار انتخاب نمایم. عمده کار این پایان‌نامه صرف نگارش برنامه اصلی پروژه گردید. روشی که در برنامه اتخاذ شده است، عمدتاً متکی است بر روش معرفی شده در شوتس (۱۹۹۵) ولی برای این که با محدودیت‌های انجام این پروژه سازگار شود، تعدیلهایی در روش مزبور اعمال شده است که طبیعتاً از دقت بخش خودکار سیستم می‌کاهد. بخش دوم کار مربوط بود به تدوین یک مجموعه برچسب که با استفاده از آن بتوان متن را برچسب‌دهی نمود. مجموعه برچسب باید به گونه‌ای طراحی می‌شد که باروش به کار رفته در پروژه سازگار باشد و نیز از معیارهای عمومی برای ایجاد چنین مجموعه‌ای تخطی نکند. بخش سوم کار، به آزمایش پروژه و بررسی نتایج اختصاص یافت.

پایان‌نامه حاضر شامل پنج فصل است. فصل اول به معرفی رشته زبانشناسی پیکره بنیاد می‌پردازد که برچسب‌دهی دستوری شاخه‌ای از آن است. در این فصل گفته می‌شود که در زبانشناسی پیکره بنیاد یکی از اعمال اصلی که روی پیکره انجام می‌شود، نشانه‌گذاری و آماده کردن آن برای کاربردهای گوناگونی است که می‌تواند داشته باشد. در فصل دوم انواع نشانه‌گذاری معرفی می‌گردد. فصل سوم اختصاص دارد به برچسب‌دهی دستوری که نوع خاصی از نشانه‌گذاری پیکره و موضوع اصلی این پروژه است. در فصل مذکور روش اصلی در این کار معرفی می‌گردد و روش اتخاذ شده شرح داده می‌شود. فصل چهارم که بخش عمده‌ای از پایان‌نامه را به خود اختصاص داده است، شامل معرفی پروژه برچسب‌دهی دستوری فارسی است. در بخش اول آن معیارهای ایجاد یک مجموعه برچسب و سپس مجموعه برچسب طراحی شده برای پروژه ارائه می‌گردند و در بخش دوم آن برنامه‌های نوشته شده و برنامه‌های به کار رفته در پروژه تک تک معرفی و شرح داده می‌شوند. فصل پنجم، شرح آزمایش‌های انجام شده روی بسته نرم‌افزار و بررسی نتایج به دست آمده را شامل

شده و در انتهای فصل مزبور پیشنهادهایی برای پروژه‌های بعدی مطرح می‌شود. در اینجا جا دارد از کسانی که نگارنده را در به انجام رسانیدن این پروژه یاری کرده‌اند قدردانی شود. ابتدا از جناب آقای دکتر سید مصطفی عاصی تشکر می‌کنم که با گشاده‌رویی و با شکیبایی فراوان مزاحمت‌های مکرر مرا تحمل کرده و از ارائه هیچ‌گونه کمکی دریغ ننموده‌اند. زیانشناسی تنها گوشه‌ای از مطالبی است که بنده از ایشان فراگرفته‌ام. دیگر از جناب آقای دکتر شاپور اعتماد متشکرم که علیرغم گرفتاری‌های فراوان و ضیق وقت مشاوره‌کار را برعهده گرفتند. از برادرم آقای علی حاجی عبدالحسینی بی‌نهایت سپاسگزارم که دستگاه رایانه خود را در اختیارم قرار داده و مزاحمت‌های ممتد مرا تا پاسی از شب گذشته متحمل می‌شدند. از پرسنل مرکز کامپیوتر پژوهشگاه علوم انسانی و مطالعات فرهنگی به دلیل همکاری‌هایشان ضمن جستجو برای منابع و نگارش، غلط‌گیری و چاپ برنامه‌ها متشکرم. در پایان از سرکار خانم محجوب که زحمت تایپ پایان‌نامه را متقبل شدند و سرکار خانم صفوی که امکان چاپ رایانه‌ای آن را برای نگارنده فراهم ساختند بی‌نهایت سپاسگزارم.

محمد حاجی عبدالحسینی

۱۳۷۵/۱۱/۲۴

## فصل ۱ زبانشناسی پیکره بنیاد

### ۱-۱ درآمد

زبانشناسان همواره بر بنیاد نهادن بررسی‌های زبانی و استدلال‌ها و استنتاج‌های مبتنی بر داده‌های واقعی زبان تأکید داشته و دارند، و همواره بهترین راه برای استفاده از داده‌های مستند را ایجاد و استفاده از پیکره‌های زبانی می‌دانسته‌اند. از طرفی گسترش روزافزون رایانه و علوم مختلف میان رشته‌ای، این امکان را برای زبانشناس فراهم می‌آورد که بتواند پیکره‌های بسیار عظیم زبانی را بوجود آورده، آنها را به گونه‌های مختلف طبقه‌بندی و نشانه‌گذاری نموده و به طرق بسیار گوناگون از اطلاعات به دست آمده بهره‌مند گردد.

گنجایش سیستم‌های ذخیره‌سازی اطلاعات و سرعت پردازش رایانه‌های امروزی به حدی بالا رفته است که به پژوهشگر امکان می‌دهد تا کاری را که گاهی ممکن بود ماه‌ها به طول انجامد و مخارج زیادی در بر داشته باشد، ظرف مدت چند دقیقه و با دقت بالا به انجام رساند. به همین دلیل، در طی دهه‌های ۱۹۸۰ و ۱۹۹۰ شاخه‌ای جدید از زبانشناسی شروع به رشد کرد که امروزه آن را با نام «زبانشناسی پیکره بنیاد»<sup>۱</sup> می‌شناسیم. «گرایش فعلی در زبانشناسی پیکره بنیاد به سوی ایجاد پیکره‌های متنی هر چه بزرگتر است» (پیچی<sup>۲</sup>، ۱۹۹۴: ۵۰۱). هلمس - هیگین<sup>۳</sup> و دیگران (۱۹۹۴: ۳۹۰) اظهار می‌دارند که «استفاده از پیکره‌های متنی، بویژه استفاده از پیکره‌های متنی رایانه‌ای، کاربردهای مفید ویژه‌ای برای مطالعه زبان‌ها و شاید در حدی پایین‌تر برای آموزش و یادگیری زبان داشته است.» از میان کاربردهای متعدد پیکره‌های زبانی می‌توان به اشاره به موارد زیر بسنده کرد: پیکره‌ها ممکن است؛

۱- منبعی باشند برای پژوهش‌های زبانشناختی و آموزش زبان،

۲- منبعی باشند برای فرهنگ‌نگاری،

۳- در سیستم‌های درک گفتار طبیعی به کار روند،

۴- در سیستم‌های ترجمه براساس اطلاعات آماری به کار روند، و یا

1. Corpus Linguistics

2. Picchi

3. Holmes-Higgin



۵. منبعی باشند برای ارزیابی و استنتاج خودکار یا نیمه خودکار دستورهای زبان (گارساید<sup>۴</sup>، ۱۹۹۳: ۳۹).

واضح است که با وجود پیکره‌های زبانی عظیم و در حال رشد و با وجود کاربردهای بسیار گوناگون برای این پیکره‌ها، نمی‌توان آنها را به شکل خام و مانند مجموعه‌ای ساده از متون مختلف نگهداری نمود، بلکه این متون می‌باید دسته‌بندی و به روش‌های مختلف نشانه‌گذاری گردند تا بتوانند پاسخگوی نیازهای کاربران باشند. در واقع ایجاد یک پیکرهٔ متنی رایانه‌ای یعنی تبدیل متون منتشر شده به یک بانک اطلاعاتی. هلمس - هیگین و دیگران (۱۹۹۴: ۳۹۰) مراحل زیر را برای این تبدیل می‌شناسند:

۱- رمزگذاری<sup>۵</sup>

۲- توصیف<sup>۶</sup> متن

۳- و در صورت امکان بازنمایی<sup>۷</sup> متن

رمزگذاری الکترونیک متن، در اصل عبارت است از علامت‌گذاری قسمت‌های متفاوت متن از جمله اطلاعات مربوط به پیکربندی آن و غیره، به‌طوری‌که بتوان اطلاعات مربوط به پیکربندی را از خود متن جدا نمود. در بخش ۱-۲-۱ در این مورد بیشتر توضیح داده خواهد شد. توصیف متن به این معنی است که محتوای متن‌های درون یک پیکره برای کاربران توضیح داده می‌شود. این کار خود نیزمند طبقه‌بندی متون موجود و ایجاد امکان دسترسی به متون براساس این طبقه‌بندی (ها) برای کاربران از طریق یک برنامهٔ مدیریت پیکره<sup>۸</sup> است.

مرحلهٔ سوم، یعنی بازنمایی متن، از دو مرحلهٔ پیشین بسیار پیچیده‌تر است و طی آن باید قراردادهای نحوی و معنایی خاصی را مشخص کرد که مطابق آنها بتوان متن را روی یک سیستم رایانه‌ای بازنمایی کرد. زمانی که یک متن در چنین سیستمی بازنموده شد، امکان این وجود خواهد داشت که یک نرم‌افزار دیگر بتواند اصلاحات تازه‌ای را از متن استنتاج کند؛ در واقع به تعبیری می‌توان گفت که رایانه خواهد توانست با استفاده

4. Garside

5. coding

6. description

7. representation

8. corpus management program

این قرارداد ادعای نحوی و معنایی متن را درک کند (هلمس - هیگین و دیگران، ۱۹۹۴: ۳۴-۳۵). برای رمزگذاری و بازنمایی متن باید آن را طبق قراردادهای خاصی شناسه‌گذاری نمود.

### ۶-۱ نشانه‌گذاری پیکره

جفری لیچ<sup>۱۱</sup>، استاد زبانشناسی و زبان انگلیسی نوین در دانشگاه لنکستر<sup>۱۱</sup> بریتانیا، در مقاله ۱۹۹۳ خود (ص ۲۷۵) نشانه‌گذاری پیکره را این چنین تعریف می‌کند: «نشانه‌گذاری پیکره عبارت است از عمل افزودن اطلاعات تفسیری (بویژه اطلاعات زبانی) به یک پیکره گفتاری و / یا نوشتاری زبان، با استفاده از نوعی رمزگذاری بر روی نمود<sup>۱۲</sup> الکترونیکی ماده زبانی».

ز آنجکه اطلاعات تفسیری می‌توانند انواع مختلفی داشته باشند، پس انواع مختلف نشانه‌گذاری بر وجهی درت (لیچ، ۱۹۹۳: ۲۷۵-۲۷۸):

#### ۱- نگارشی<sup>۱۳</sup>

#### ۲- آوایی و جوی

#### ۳- بیرونی

#### ۴- دستوری (بر حسب دهم دستوری)<sup>۱۴</sup>

#### ۵- نحوی تنظیم<sup>۱۵</sup>

#### ۶- معنایی

#### ۷- بازکردن ساختار کلامی

### ۶-۱ نشانه‌گذاری نگارشی

نوعی از نشانه‌گذاری که هیچ معرفی می‌کند یعنی نشانه‌گذاری نگارشی در واقع همانند رمزگذاری متن است که هلمس - هیگین و دیگران به آن اشاره کرده‌اند (ر.ک. ۱۱). برای نشانه‌گذاری نگارشی ابتدا نیاز به یک زبان علامت‌گذاری<sup>۱۶</sup> وجود دارد. بر مبنای چنین زبانی جزء یک متن بر اساس نقش خود علامت‌گذاری می‌شوند. اجزاء متن می‌توانند هر عنصری از متن باشند مانند یک پاراگراف، چکیده، یادداشت و یا

9. annotation

10. Geoffrey Leech

11. Lancaster

12. representation

13. orthographic

14. grammatical taggery

15. parsing

16. markup language

شکل ۱-۱ نمونه‌ای از علامت‌گذاری SGML (اسمیت، ۱۹۸۷: ۱۷۲)

```
<sc
<mb
<h1 The Standard Generalized Markup Language (SGML) for
Humanities Publishing
<au Joan M Smith
<ad
<l>National Computing Centre
</al>
<ab
<p>A new methodology, at the core of which is generic coding,
has been developed within the International Organization for
Standardization (ISO).
This is known as the &SGML;.
Using SGML, the elements of a document are marked up as to
their role, be it a paragraph, an abstract, a note, or
whatever; the style of presentation is a separate issue and
is not addressed by SGML.
These elements can form part of a data base, which can be
updated as well.
So there is the notion of data base publishing.
<p>The &SGML is presented as a tool for full&en;text data
base publishing, where the options for output are open, an
example being given of a marked up document.
Its value for all aspects of humanities publishing is
addressed: whether for scholarly papers intended for a
journal, books, specialist publications, dictionaries, or
bibliographies; indeed, whatever is input to an electronic
medium with the intention of its being imaged subsequently in
some form- whether alone or in combination with other text.
SGML represents an advance in publishing methodology, taking
advantage of developing technology.
It can be exploited as such in an academic environment to
give an added dimension to research publications.
```

هر عنصر دیگر (اسمیت<sup>۱۷</sup>، ۱۹۷۸، ص ۱۰۰). نگارگری که حتی - استفاده از یک زبان علامت‌گذاری، نشانه‌گذاری نگارگری گردید. نکته‌ای که در آن متن را بسادگی به رایانه‌های دیگر منتقل نمود.

از بین زبان‌های علامت‌گذاری می‌توان به SGML<sup>۱۸</sup> که توسط سازمان استاندارد جهانی (ISO) ساخته شده و ISO ۳۰۱۱<sup>۱۹</sup> نامیده می‌شود و رایج است بین‌المللی که تحت حمایت مالی انجمن رایانه عمومی<sup>۲۰</sup>، انجمن زبانشناسی رایانه‌ای<sup>۲۱</sup> و انجمن زبانشناسی و ادبیات رایانه‌ای<sup>۲۲</sup> فعالیت دارند.

شکل ۱-۱ نمونه‌ای از یک متن انگلیسی را که - SGML<sup>۲۳</sup> شده‌گذاری شده نشان می‌دهد. همان‌طور که ملاحظه می‌شود نشانه‌های < > قرار داده می‌شوند تا رایانه بسادگی بتواند آن را از خود متن تشخیص دهد. به عنوان مثال در نمونه مزبور <sd> به معنی شروع شدن (start document) است و <mb> یعنی بدنه اصلی متن (main body).

بدین ترتیب رایانه می‌تواند با هر یک از جزء متن به شکل خصوصی برخورد نماید، مثلاً بخشی را با قلم خاصی بنویسد و یا فقط قسمت‌های خاصی از متن را نمایش داده یا به رایانه دیگری منتقل نماید.

#### ۲-۲-۱ نشانه‌گذاری آوایی / واجی

این نوع نشانه‌گذاری در علوم تشریحی و تکثیرایی گسترده‌ای دارد.

#### ۳-۲-۱ نشانه‌گذاری زیر زنجیری

لیچ (۱۹۹۳: ۲۷۶) اظهار می‌دارد که تریکرت<sup>۲۴</sup> شده‌گذاری شده از لحاظ ویژگی‌های زیر زنجیری وجود دارد: ۱- پیکیته شدن - لاند<sup>۲۵</sup> که تنها بیست‌هزار واژه انگلیسی

17. Joan M. Smith

18. Standard Generalized Markup Language

19. Text Encoding Initiative

20. The Association for Computing and Humanities

21. The Association for Computational Linguistics

22. The Association for Literary and Linguistic Computing

23. The London-Land Corpus

بریتانیا می‌شود، ۲- پیکره انگلیسی گفتاری لنکستر/آی بی ام<sup>۲۴</sup>؛ این پیکره تنها شامل پنجاه هزار واژه است اما در سطوح مختلف زبانی نشانه گذاری شده است. این گونه نشانه گذاری دستی انجام می‌شود و مهارت نشانه گذار در تشخیص ویژگی‌های زیرزنجیری درست، از اهمیت خاصی برخوردار است. شکل ۲-۱ نمونه‌ای از یک مکالمه نشانه گذاری شده را به همراه معانی نشانه‌ها نمایش می‌دهد.

#### ۴-۳-۱۱ برجسب دهی دستوری

این نوع نشانه گذاری که موضوع همین پایان‌نامه است در فصل‌های بعدی به تفصیل توضیح داده خواهد شد.

#### ۵-۳-۱۱ تقطیع نحوی

غالب به برجسب‌های دستوری به عنوان مرحله اول یک نشانه گذاری نحوی جامع تر نگریسته می‌شود. این نشانه گذاری عبارت است از تهیه یک نمایش درختی<sup>۲۵</sup> یا قلاب‌های نشاندار<sup>۲۶</sup> برای هر یک از جمله‌های درون پیکره. پیکره‌های تقطیع شده نحوی را با یک درختی<sup>۲۷</sup> می‌نامند. بانک‌های درختی گوناگون و کوچکی تهیه شده‌اند و به منظور افزایش سرعت و صحت نشانه گذاری گرایش به این سو پیدا شده است که از سطح‌های ساده‌تری استفاده شود. شکل‌های ۳-۱ و ۴-۱ دو نمونه از تقطیع جملات نگلیسی را نشان می‌دهند.

#### ۶-۳-۱۱ نشانه گذاری معنایی

این نشانه گذاری نحوی (یعنی برجسب دهی دستوری و تقطیع جملات) روشن است که قدم بعدی نشانه گذاری معنایی است. برجسب دهی معنایی کلمات مثلاً<sup>۲۸</sup> می‌توانند به این هدف انجام شود که معنای واژگانی کلمات در متن شناسایی شوند؛ این روشی است که به آن تعیین معنی<sup>۲۸</sup> می‌گویند. شکل ۵-۱ نمونه‌ای است از برجسب‌دهی معنایی کلمات که از یک پروژه تحلیل خودکار معنا در دانشگاه لنکستر گرفته شده است. در این مثال، متن از بالا به پایین خوانده می‌شود: برجسب‌های دستوری در سمت چپ قرار دارند و برجسب‌های معنایی در سمت راست.

24. The Lancaster/IBM Spoken English Corpus

25. phrase marker

26. labelled brackets

27. treebank

28. sense resolution

شکل ۲-۱ نشانه گذاری زیر زنجیری، پیکره لندن - لاند

توجه: البته این نسخه ASCII می تواند برای کاربر نهی (تیس) (تیس) (لیج) (۱۹۹۳: ۲۷۶)

```
well ^very nice of you to ((come and)) _spare the
^me and#
^come and !(\alk# -
^tell me about the - !pr\oblems#
and ^inci\dentially# .
^do ^do !\ell me#
^anything you 'want about the :college in "g\eneral#
^mean it 'doesnt "h\ave to be (con^fined#)#
to the ^problems of !\English#
^er and the ^horrors of :living in :this 'part
of the c/ollege#
or ^anything like th\at#
* - laugh)*
```

KEY

- are features of stress, including boosters
- ^ are tones (fall, rise, and fall-rise)
- \_ are pauses
- # is a tone unit boundary
- (B) is a pause filler 'er'
- { } enclose a subordinate tone unit
- [ ] enclose contextual comments
- [ ] enclose uncertain material

شکل ۳-۱ تقطیع پیرایشی (Skeleton parsing)، از یک پیکره گفتاری زیر زنجیری

فلاهای بی نام نشان دهنده سازه ای هستند که در طرح نشانه گذاری زیر زنجیری ارائه نشده است. (نشانه گذارها مجاز هستند که بدون نشان دادن محتوی ک سازه ها به آنها تعلق دارند، آنها را مشخص نمایند.) (لیج، ۱۹۹۳: ۲۷۸)

```
[S[N the_AT killer_NN1 whale_NN1 N] _ [Fr[N
whc_PNQS ^ [V 'd_VHD grown_VVN ] ] too_RG big_I] [P for_IF
[N re_APP5 pool_NN1 [P on_I] [N Clacton_NP1 Pier_NNL1
NPNP]]] [V has_VHZ arrived_VVN safely_RR [P at_I
[N re_APP5 new_I] home_NN1 [P in_I] [N Windsor_NP1 | safari_NN1
pan_NNL | NPNP] ] _ S]
```

KEY: Square brackets enclose constituents above word level. Brackets are linked to their words by '-'. The tagset used here is the revised version of the LOB tagset, known as the 'CLAWS2 tagset'. Tag definitions can be inferred from a comparison between this example and Fig. 2. Symbols for constituents above word level are: Fr: relative clause; J: adjective phrase; N: noun phrase; P: prepositional phrase; S: sentence; V: verb phrase

شکل ۴-۱. تصحیح پیرینسی، زانک و جونیور (Penny, Zank & Juniors: ۱۹۹۳: ۲۷۸)

(۱۳)

NP M. linker  
 NP is  
 NP chairman  
 PP c  
 NP (NP Elsevier ...)

NP the Dutch  
 publishing  
 group) ( )

شکل ۵-۱. نوعی برجسته شعر معاصر از احمد نوری (۱۹۹۳: ۲۷۱)

|       |              |            |
|-------|--------------|------------|
| NP    | the          | Z8         |
| AVE   | laughed      | E4.1+      |
| RF    | disagreement | O4.2-      |
| AVE   | quashing     | A1.1.1     |
| APP2E | ter          | Z8         |
| NN    | cigarette    | F3         |
| I     | it           | Z5         |
| AT    | he           | Z5         |
| NN    | litter       | F1         |
| CC    | and          | Z5         |
| AVE   | vent         | M1         |
| II    | o            | Z5         |
| AT    | he           | Z5         |
| NN    | telephone    | Q1.3       |
| TO    | o            | Z5         |
| AV    | ring         | Q1.3[1.2.1 |
| RF    | o            | Q1.3[1.2.2 |
| DE    | al           | N5.1+      |
| APP2E | ter          | Z8         |
| NN    | needs        | S2/S3.1    |

KEY: A1.1.1, general actions; E4.1, happy and sad; F1, food; F3, cigarettes and drugs; M1, motion; N5.1, all, none; O4.2, judgement of appearance; Q1.3, telecommunications; S2, people; S3.1, relationships, general; Z5, function words; Z8, pronouns, etc.

شکل ۶.۱ نشانه گذار کلامی و ربط نحوی در خبرهای سوئیته پرس (Associated Press)

(لیج، ۱۹۹۳: ۲۰۹)

S.1 (0) The state Supreme Court has refused to release [1 [2 Rahway State Prison 2] inmate 1] [1 James Scott 1] on bail.

S.2 (1) The fighter [1 is serving 30-40 years for a 1975 armed robbery conviction

S.3 (1) Scott [1] has asked for freedom while <1 he waits for an appeal decision.

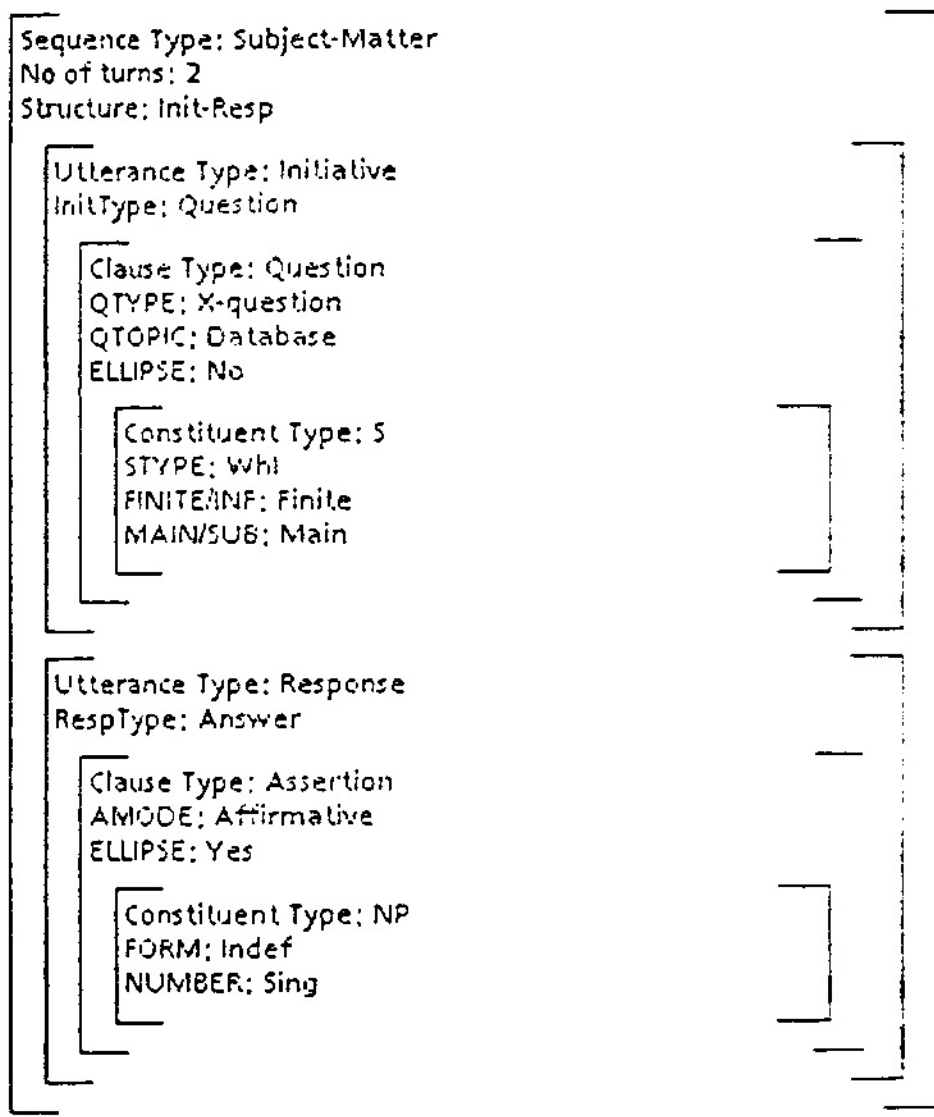
S.4 Meanwhile [1 <1 his promoter 3] [3 Murad Muhammed 3], said Wednesday <3 he netted only \$13,250 for [4 [1 Scott 1's nationally televised light heavyweight fight against [3 ranking contender 5]] (5 Yaqui Lopez 5) last Saturday 4).

S.5 (4) The fight in which [1 Scott 1] won a unanimous decision over (5 Lopez 5) [4] grossed \$135,000 for [3 Muhammed 3]'s firm 6], [6 Triange Productions of Newark 6], [6] he said.

KEY: The use of the same index (i, 2, ...) binds one syntactic constituent to another to which it is co-referential or semantically equivalent. In the following list, [ ] represents an arbitrary index:

- (i i) or [i i] enclose a constituent (normally a noun phrase) entering into an equivalence chain
- <i indicates a pronoun with a preceding antecedent
- >i indicates a pronoun with a following antecedent
- {(i i)} enclose a noun phrase entering into a copular relationship with a preceding noun phrase
- {i i)} enclose a noun phrase entering into a copular relationship with a following noun phrase
- (0) represents an anaphoric barrier: in effect, the beginning of a new text.





## ۷-۲-۱ نشانه گذاری کارکردشناختی / کلامی

اصطلاح پوششی «نشانه گذاری کلامی» را می توان برای هر نوع نشانه گذاری که به واحدها یا روابطی وزای مرز جمله برمی گردد، به کار برد. برای داده هایی از زبان گفتاری نوعی نشانه گذاری به کار رفته است (اشتن اشتروم<sup>۲۹</sup>، ۱۹۸۴) تا طبقاتی از علائم کلامی را که نشان دهنده ساختار مکالمه می باشند نشان دهد. نوع دیگری از نشانه گذاری که تا زمان نگارش مقاله (لیچ، ۱۹۹۳) تنها روی داده های نوشتاری در دانشگاه لنگستر به کار گرفته شده بود، نشان دهنده روابط ارجاعی<sup>۳۰</sup> و روابط همسنگی<sup>۳۱</sup> دیگر میان اجزاء همسایه در یک متن است. نمونه ای از متنی که به این روش نشانه گذاری شده در شکل ۶-۱ آمده است. آهرن برگ و یونسن<sup>۳۲</sup> (۱۹۸۸: ۶۶-۷۰) از دانشگاه لینکوپینگ<sup>۳۳</sup> سوئد نیز یک سیستم محاوره ای<sup>۳۴</sup> را برای برجسب دهی مکالمات معرفی می کنند که در آن برجسب ها برای نشان دادن اطلاعات کارکرد شناختی به کار می روند. این اطلاعات برای مثال عبارتند از: نشان دادن برنهاده<sup>۳۵</sup> (ها) در یک نقطه خاص از کلام یا کنش (های) گفتاری که در یک پاره گفتار ویژه اجرا می شود. شکل ۷-۱ نمونه ای از ساختار برجسب های این سیستم را که Dagtag نام دارد، نشان می دهد.

### ۳-۱ چکیده

در این فصل شاخه میان رشته ای زبانشناسی پیکره بنیاد به طور اجمالی معرفی گردید و ضمن ذکر مراحل ایجاد یک پیکره رایانه ای، مسأله نشانه گذاری پیکره نیز به میان آورده شد. گفته شد که نشانه گذاری پیکره عبارت است از افزودن اطلاعات تفسیری عمدتاً زبانی به یک پیکره گفتاری و یا نوشتاری زبان با استفاده از نوعی رمزگذاری بر روی نمود الکترونیک ماده زبانی. سپس انواع نشانه گذاری پیکره به اختصار معرفی گردید. فصل بعدی اختصاص دارد به معرفی برجسب دهی دستوری پیکره.

29. Stenström

30. anaphoric relations

31. equivalence relations

32. L.Ahrenberg and A.Jönsson

33. Linköping

34. interactive

35. topic

## فصل ۲ برچسب دهی دستوری

برچسب دهی دستوری، یا مشخص نمودن مقوله دستوری هر واژه، نوعی نشانه گذاری است که بیش از همه در پیکره های تهیه شده در اروپا و آمریکا رایج است و این خود دو دلیل دارد:

۱- این کار به اندازه ای ساده هست که بتوان آن را تا حد زیادی به صورت خودکار انجام داد.

۲- کاربردهای فراوانی دارد، مثلاً یکی در فرهنگ نگاری که در آن برچسب دهی دستوری قدم اول است در عمل سرواژه یابی<sup>۱</sup> (لیچ، ۱۹۹۳: ۲۷۶). ضمناً همان گونه که در بخش ۱-۲-۵ ذکر شد، این کار اولین گام در تقطیع نحوی جملات نیز هست.

یک طرح برچسب دهی دستوری شامل موارد زیر است (لیچ، ۱۹۹۳: ۲۷۶):

۱- یک مجموعه برچسب<sup>۲</sup>: مجموعه ای از عنوان های دستوری برای اقسام واژگانی

۲- مجموعه ای از تعاریف برچسب ها

۳- مجموعه ای از دستورالعمل های برچسب دهی که نحوه افزودن برچسب ها را به متن مشخص می کند.

هنگامی که اقسام واژگانی تحلیل شدند و مقوله دستوری هر یک مشخص شد، برچسب مربوط به مقوله دستوری هر واژه - مانند هر نوع عمل نشانه گذاری دیگر (ر.ک. ۱-۲) - در کنار آن واژه وارد متن می گردد. البته این کار باید به طریقی انجام شود که باز بتوان بدگی برچسب ها را حذف کرد و به متن اصلی دست یافت. شکل ۱-۲ نمونه ای از یک متن برچسب داده شده به همراه معانی برچسب های آن را نشان می دهد.

### ۱-۲ تاریخچه نشانه گذاری نحوی

عمل برچسب دهی دستوری یک پیکره زبانی تاکنون در ایران انجام نپذیرفته است. حتی تهیه یک پیکره نیز از سابقه ای طولانی در کشور ما برخوردار نیست، بلکه از اوایل

شکل ۱-۲ برچسب‌دهی دستوری، پیکره LOB توجه کنید که در هر طرح، معنی طرح‌های دیگر برای متون نوشتاری، علائم نقطه‌گذاری (وزنه) به حساب آید و برچسب‌دهی می‌شوند.

maintain\_NN is\_BEZ an\_AT excellent\_JJ virtue\_NN  
 but\_CC not\_XNCT when\_WRB the\_ATI guests\_NNS have\_HV  
 a\_TO seat\_VB n\_IN rows\_NNS in\_IN the\_ATI cellar\_NN

he\_AT covers\_NNS whose\_WP\$ chief\_JJB scene\_NN  
 was\_BEZ cut\_VEN at\_IN the\_ATI last\_AP moment\_NN  
 had\_HVD comparatively\_RB little\_AP to\_TO sing\_VB

she\_PP3A stole\_VBD my\_PPS wallet\_NN  
 named\_VBD Rollinson\_NP

Example: AT post-determiner; AT, singular article; ATI, article neutral or number; BEZ, was (past sing. form of the verb *be*); BEZ, is (3rd form of the verb *be*); CC, coordinating conjunction; HV, base form of the verb *have*; HVD, had, 'd (past form of the verb *have*); IN preposition; JJ, general adjective; JJB, attributive adjective; NN, singular common noun; NNS, plural common noun; NP, singular proper noun; PP3A, third-person singular personal pronoun, subjective; PPS, possessive pronoun (determinative function); RB, general adverb; TO, to-infinitive marker; VB, base form of lexical verb; VBD, past tense of lexical verb; VEN, past participle of lexical verb; WPS, possessive *wh*-pronoun *whose*; XNCT, negative particle (*not, n't*)