

نرم افزار تشخیص فعل در زبان فارسی *

پویان دانش کار آراسته
کارشناس ارشد زبان شناسی

درآمد

زبان شناسی رایانه‌ای^۱ از پویاترین شاخه‌های علم زبان شناسی در جهان است که ناشناخته‌تر از دیگر حوزه‌های این علم در ایران است. این مقاله، گزارشی است از تلاشی عملی که در این زمینه به انجام رسیده است. در این مقاله، روشی در طراحی نرم‌افزاری ارائه می‌شود که می‌تواند با کمک مشخصات صوری فعل در زبان فارسی، این مقوله را در متن تشخیص دهد، تقطیع تکوازی کند، و مشخصات صرفی و نحوی آن از قبیل زمان، شخص، عدد، مثبت و منفی، معلوم و مجهول، سببی، وجهی، امر، و مرکب را اعلام کند.

* این مقاله براساس رساله کارشناسی ارشد نگارنده - با همین عنوان - به راهنمایی و مشاوره استادان ارجمند جناب آقای دکتر سید مصطفی عاصی و جناب آقای دکتر محمد دبیرمقدم که در بهمن ۱۳۸۱ در دانشگاه علامه طباطبایی دفاع شد، به نگارش درآمده است. علاقه‌مندان به مشاهده و تحقیق درباره کُد برنامه که به زبان Visual Basic نوشته شده است، می‌توانند به این رساله مراجعه فرمایند. در ضمن، از آقای ایمان شاکریان که در گُذونویسی این برنامه مرا بسیار یاری کردند، صمیمانه تشکر می‌کنم.

1. Computational Linguistics

فرهنگ، ۵۰-۴۹؛ بهار و تابستان ۱۳۸۳؛ ص ۴۶-۳۱

در این مقاله فقط جنبه زبان‌شناختی مسئله بررسی می‌شود. جنبه‌های تخصصی رایانه‌ای کار، مستلزم مجالی دیگر و مقاله‌ای مستقل است. این طرح می‌تواند برای نرم‌افزارهای مشابه، به‌منزله یک تجربه انجام‌شده، راهگشا باشد.

۱. مقدمه

بی‌شک یکی از اهداف زبان‌شناسی رایانه‌ای که از ابتدای تکوین و گسترش این علم مد نظر بوده و پیوسته دنبال شده، ارائه نرم‌افزارهایی است که به کمک آنها بتوان تحلیل‌هایی بر روی متن و یا پاره‌های زبانی انجام داد. تولید این‌گونه نرم‌افزارها در بلندپروازانه‌ترین شکل خود، برنامه‌های جامع ترجمه ماشینی و برنامه‌های تحلیل متن است.

روشن است که یک برنامه جامع، مثل ترجمه ماشینی یک زبان به زبان دیگر، برنامه‌ای فوق‌العاده گسترده است که خود از اجتماع چندین برنامه کوچک‌تر که هر یک وظیفه‌ای خاص برعهده دارد، تشکیل می‌شود. از عمده‌ترین اجزای چنین برنامه‌هایی، نرم‌افزارهایی است که بتواند مقولات زبانی را تشخیص و مشخصات آنها را به‌طور دقیق ارائه دهد تا برنامه اصلی، بنابه هدفی که دارد، آنها را تحلیل و یا ترجمه کند.

یکی از اصلی‌ترین و درعین حال پیچیده‌ترین مقولات زبان فارسی، فعل است؛ که در مقاله حاضر درباره آن بحث می‌شود.

قطعاً اجرای چنین تحقیقاتی، علاوه بر طبیعت کاربردی، به دستاوردهای نظری و احیاناً شبیه‌سازی و درک فرایندهای زبانی مقولات مورد نظر نیز می‌انجامد.

در این مقاله ابتدا به صورتی بسیار گذرا، از مقوله فعل و آنچه این نرم‌افزار در مورد آن باید تشخیص دهد سخن به میان می‌آید و بحثی کوتاه درباره افعال مرکب ارائه می‌شود؛ سپس ساختار کلی برنامه معرفی و تک‌تک مراحل آن بررسی می‌شود؛ در پایان نیز ویژگی‌ها و مشکلات کار به اجمال مطرح و نتیجه‌گیری می‌شود.

۲. فعل در زبان فارسی

مقوله فعل در زبان فارسی بیش از هر مقوله دیگر، مشخصات صوری دارد: وندهای متعدد، از قبیل پیشوندها، شناسه، تکواژ سببی ساز، و مانند آن. این امر از یک سو نوشتن یک برنامه رایانه ای را که طبعاً به عوامل صوری و قابل تعریف برای رایانه نیاز دارد، آسان و از سوی دیگر، به دلیل تعدد این مشخصات و گستردگی و حالات بسیار زیاد فعل در زبان فارسی، فوق العاده پیچیده می کند.

نرم افزار مورد نظر باید صورت های مختلف فعلی در زبان فارسی را که در زیر فهرست وار بیان شده است، تشخیص دهد:

الف- تمامی گروه فعلی زبان فارسی، نه فقط فعل واژگانی؛ برای مثال: می خوانند بروند، گفته نشده بود، داشتیم می خوردیم، نپوشانده بودی، فرو افتاده بود، و...

ب- صورت خاص هر یک از زمان های فعل فارسی از قبیل حال ساده، حال استمراری، ماضی ساده، ماضی بعید، آینده، ...

بنابراین تمامی زمان های فعل فارسی باید تجزیه و تحلیل شود و مشخصات صوری هر زمان به شکل فرمول های قابل تعریف برای رایانه درآید. برای نمونه، فرمول زمان حال ساده و صورت سببی مجهول زمان گذشته استمراری به ترتیب چنین بیان می شود:

می + بن حال + شناسه (م، ی، د، یم، ید، ند)

داشت + شناسه (م، ی، یم، ید، ند) + بن حال + انده/انیده + می + شد + شناسه (م، ی، یم، ید، ند)

(یم، ید، ند)

یادآور می شود که چون تحلیل فعل فارسی در متن مورد نظر است، باید صورت نوشتاری هر کلمه و یا تکواژ در نظر گرفته شود.

پ- صورت مثبت و منفی هر صورت فعلی؛ یعنی تکواژ منفی ساز و تمامی جایگاه های قرارگرفتن آن باید برای رایانه تعریف شود. برای مثال: نرفتیم، نخواهد خواند، گفته نخواهد شد، گفته نشد، و... (چنان که ملاحظه می شود، تکواژ منفی ساز در زبان فارسی جایگاه های مختلفی دارد که هر یک باید مشخص شود).

ت- صورت مجهول افعال؛ که باید به صورت یکپارچه تشخیص داده شود، یعنی مثلاً فعل «خورده خواهد شد» به طور کلی تشخیص داده شود، نه به صورت جدا جدا، مثلاً «خواهد» و «شد».

ث- فعل‌های سببی، که تکرار سببی ساز دارند؛ مانند: بخورانیم، خورانیده شد... همان‌گونه که در قسمت «ب» عنوان شد، بنابر ملاحظات رایانه‌ای، صورت نوشتاری عنصر سببی ساز به سه صورت ان/انده/انیده تعریف می‌شود.

ج- صورت‌های امر و نهی افعال؛ مثل: برو، نرو، بسوزان...

چ- صورت‌های نوشتاری مثبت و منفی در صیغه سوم شخص مفرد زمان‌های زیر در افعالی که بن ماضی آنها در نوشتار، از اضافه کردن یک «د» به بن حال به دست می‌آید، مثل فعل «خوردن»، یکسان است، که باید هر دو زمان منظور شود:

حال ساده (/سببی) = ماضی استمراری (/سببی)

حال التزامی (/سببی) = ماضی ساده (/سببی)

در ضمن، ظاهر مصدر کوتاه با صیغه سوم شخص مفرد فعل ماضی ساده یکی است.

ح- افعال وجهی مانند باید (می‌باید، بایستی، می‌بایستی، می‌بایست)، می‌شود، بشود، می‌شد، می‌توان، بتوان، همچنین صورت‌های ماضی ساده، ماضی استمراری، ماضی نقلی، ماضی التزامی، و نیز حال ساده از فعل توانستن (مشکوة الدینی، ۱۳۷۴: ۱۵۶).

چنان که اشاره شد، هر فعل صورت‌های زمانی مختلفی دارد. فعل را به هشت صورت مثبت، منفی، مجهول، منفی مجهول، سببی، منفی سببی، سببی مجهول، و منفی سببی می‌توان بیان کرد؛ که این صورت‌ها نیز می‌تواند برای شش شخص دستوری متفاوت باشد. در ضمن، صورت‌های امر و نهی هم وجود دارد. همچنین در بعضی از زمان‌ها ممکن است برای هر یک از موارد دو صورت وجود داشته باشد، مثل: خورانده نباشم یا نخورانده باشم.

برای مثال، صورت‌های زمان حال ساده از مصدر «خوردن» برای صیغه سوم شخص مفرد چنین است:

۱. صورت مثبت: می‌خورد

۲. صورت منفی: نمی‌خورد

۳. صورت مجهول: خورده می شود
۴. صورت منفی مجهول: خورده نمی شود
۵. صورت سببی: می خوراند
۶. صورت منفی سببی: نمی خوراند
۷. صورت سببی مجهول: خورنده/خورانیده می شود
۸. صورت منفی سببی مجهول: خورنده/خورانیده نمی شود

۱-۲. افعال مرکب

در باره افعال مرکب - که به دلیل ویژگی خاصی که دارند؛ در اینجا جداگانه درباره آنها بحث می شود - بسیاری از زبان شناسان بررسی های دقیقی انجام داده اند؛ از جمله: باطنی (۱۳۴۸)، عاصی (۱۳۷۱)، صادقی (۱۳۷۲)، دبیرمقدم (۱۳۷۴). در این زمینه، اختلافاتی وجود دارد: گروهی. همچون خانلری (۱۳۶۵: ۱۸۱-۱۱۳) و مشکوة الدینی (۱۳۷۴: ۱۶۳-۱۵۸) فعل های فارسی را به سه دسته ساده، پیشوندی، و مرکب تقسیم کرده اند؛ ولی دبیرمقدم (۱۳۷۴) فعل های پیشوندی و مرکب - هر دو - را فعل مرکب می شمارد و به تقسیم بندی فعل ساده و فعل مرکب قائل است. همچنین، خانلری (۱۳۶۵: ۱۸۱-۱۷۶) عباراتی مثل «خوشم آمد»، «خوابشان می آید» و مانند آن را تحت عنوان فعل های ناگذر تبیین می کند؛ اما به اعتقاد دبیرمقدم (۱۳۷۴)، این عبارات اساساً نه فعل مرکب بلکه جملات تمام و کمالی هستند که در آنها جزء اسمی، فاعل است.

در اینجا یادآوری یک نکته ضروری است. اصولاً در یک برنامه رایانه ای، آنچه ارزش دارد، عوامل صوری قابل تعریف برای رایانه است. از این دیدگاه، مسائل نظری تا هنگامی دخیل هستند که به کشف و تبیین این مشخصه های صوری بینجامند. بنابراین، صرف نظر از اینکه کدام یک از نظریه ها و تقسیم بندی های ذکر شده قابل قبول تر و از منظر زبان شناسی نظری توجیه پذیرتر است، در این مقاله به ترتیب زیر عمل خواهد شد:

۱. افعال مرکبی که اصطلاحاً پیشوندی نامیده می شوند، مثل «باز آمدن»، «فروافتادن»، و مانند آن، به دلیل محدود بودن پیشوندهای فعلی، قابل تعریف و

تشخیص برای رایانه هستند و در طرح تشخیص برنامه به ترتیبی که گفته خواهد شد (بخش ۴-۳-۴)، گنجانده می‌شوند.

۲. در عباراتی که اصطلاحاً ناگذر خوانده می‌شوند، تنها کلمه فعل در نظر گرفته می‌شود؛ مثال: (۱) خوابش خواهد گرفت. (خواهد گرفت: فعل آینده ساده، مثبت، صیغه سوم شخص مفرد). که البته این کار باعث افزایش درصد خطای برنامه خواهد شد؛ زیرا فعل‌های جملاتی نظیر «خوابشان خواهد گرفت»، «خوابت خواهد گرفت»، و غیره نیز همانند فعل جمله (۱) اعلام خواهد شد. چنین مواردی را به ناچار باید به پس‌ویرایش^۱ سپرد (نگاه کنید به بخش ۵ در همین مقاله).

۳. با اینکه افعالی مثل «بودن»، «کردن»، «زدن»، «آمدن»، «فرمودن»، «افتادن»، و بسیاری دیگر، در بیشتر مواقع به صورت مرکب به کار می‌روند - با توجه به بسامد بسیار بالای افعال مرکب در فارسی - مثل: «دلخور بودن»، «صبر کردن»، و مانند آن، موارد بسیاری هم وجود دارد که این افعال به صورت مرکب به کار نمی‌روند؛ ضمن اینکه بسیاری از آنها با گروه حرف اضافه‌ای که مشخص نیست از چند کلمه تشکیل می‌شود، به کار می‌روند، مانند «به دست آوردن»، «از کار افتادن». بنابراین، اگر بخواهیم کلمه یا کلمات قبل از این افعال را - که به هر حال نمی‌توان برای رایانه قابل تشخیص کرد و یا قابل تشخیص کردن آنها اقتصادی نیست - به این افعال منضم کنیم، درصد خطای برنامه بیش از حد بالا خواهد رفت. پس با توجه به اینکه برنامه مورد نظر جزء فعلی این افعال را تشخیص خواهد داد و همه اطلاعات صرفی و نحوی این افعال مانند شخص، عدد، و غیره در جزء فعلی آنها وجود دارد، بهتر است بدون آنکه درصد خطای برنامه را افزایش دهیم، تشخیص جزء غیر فعلی این افعال را به پس‌ویرایش بسپاریم.

۳. پایگاه داده

برای این برنامه، یک پایگاه داده^۲ متشکل از بن‌های ماضی و بن‌های مضارع فعل‌های زبان فارسی امروز، به صورت الفبایی، تعبیه شده است. در آینده هرگاه

نرم افزار تشخیص فعل در زبان فارسی ۳۷

کلمه جست و جوی^۱ به کار رود، منظور جست و جوی در این بانک است. نمونه‌ای از این پایگاه داده که با نرم افزار Access نوشته شده است - و می توان آن را با هر نرم افزار پایگاه داده دیگر هم نوشت - به شکل زیر است:

Presentroots	Pastroots	Number
آز	آخت	۱
آرام	آرمید	۲
آزار	آزد	۳

۴. ساختار برنامه

چون هدف در طرح مورد نظر، تشخیص فعل در متن است، مراحل کلی زیر را می توان برای گام های مختلف عملکرد نرم افزار متصور بود:

۱. گرفتن متن به عنوان درونداد

۲. تشخیص اولیه کلمه به طور مستقل

۳. بررسی فعل بودن یا نبودن کلمه

۴. اعلام مشخصات در صورت فعل بودن کلمه

۵. تشخیص و بررسی کلمه بعد

۴-۱. گرفتن متن به عنوان درونداد

در این مرحله، کاربر متنی را که ممکن است شامل یک کلمه، یک عبارت، یک جمله، و یا متنی متشکل از جمله های بسیار باشد، در محل مخصوص وارد یا کپی می کند. خوشبختانه، زبان های برنامه نویسی جدید، به ویژه برنامه هایی که از ابزار ویندوز^۲ استفاده می کنند و اصطلاحاً Visual نامیده می شوند، ابزار مورد نیاز در این قسمت مانند کادر متن^۳ را در اختیار برنامه نویس قرار می دهند.

1. Search

2. Windows

3. Text Box

۲-۴. تشخیص کلمه به صورت مستقل

ابزار تشخیص کلمه در این برنامه، فاصله^۱ است. مجموعه حروفی که بین دو فاصله قرار داشته باشد، یک کلمه در نظر گرفته می‌شود. مانند:

علی دیروز به مدرسه رفت.

گاهی در نوشتار فارسی، برخی از تکواژهای واژه فعل به صورت جدا از هم نوشته می‌شود، مانند «می‌روم»، «گفته‌اید»، و... مواردی از این دست برای برنامه تعریف شده است و برنامه به طور خودکار فاصله بین این تکواژها را حذف^۲ می‌کند و این کلمات را به صورت یک کلمه یکپارچه در نظر می‌گیرد.

۳-۴. بررسی "فعل" بودن یا نبودن کلمه

این مرحله که در حقیقت بدنه اصلی برنامه را تشکیل می‌دهد، بیشترین حجم کد برنامه را دربر می‌گیرد، و متشکل از بخش‌های^۳ متعددی است. پایه و اساس این قسمت بر این است که افعال زبان فارسی یا هیچ‌وندی - پیشوند یا پسوند - در ظاهر ندارند یا یک و یا هر دو این‌ها را دارند. بنابراین اگر بتوانیم وندهای یک فعل را جدا کنیم، آنچه باقی می‌ماند، بن مضارع یا بن ماضی خواهد بود. در این صورت می‌توانیم با شناختن وندهای به‌کاررفته در ساختمان فعل و ترتیب قرارگرفتن آنها و شناختن بن فعل، به مشخصات صرفی و نحوی فعل پی ببریم و آنها را اعلام کنیم. یک نمای بسیار کلی از این مرحله - که حدود ۳۹۰۰ خط از برنامه را دربر می‌گیرد - چنین است:

۱. بررسی ابتدایی کل کلمه

۲. جداکردن تک‌تک حروف از ابتدای کلمه و بررسی پیشوندها

۳. جداکردن تک‌تک حروف از انتهای کلمه و بررسی پسوندها

۴. عملیات میانی برای تشخیص فعل‌های چندکلمه‌ای

۴-۳-۱. بررسی ابتدایی کل کلمه

در این مرحله، کل کلمه یک بار در بانک جست و جو می شود؛ زیرا افعالی مثل «خورد» یا «رفت»، هیچ پسوند یا پیشوندی ندارند. اگر کلمه در بانک پیدا شود، قسمت اعلام مشخصات اجرا می شود.

۴-۳-۲. جدا کردن تک تک حروف از ابتدای کلمه و بررسی پیشوندها

در این مرحله، اولین حرف سمت راست - ابتدای - کلمه جدا و سپس بررسی می شود که آیا این حرف می تواند جزء اولین حروف پیشوندهای فعل فارسی باشد یا خیر؛ برای مثال، حرف «م» می تواند اولین حرف پیشوند باشد، ولی مثلاً «ج» اولین حرف هیچ پیشوند فعلی در زبان فارسی نیست. به همین ترتیب، انتخاب های بعدی نیز محدودتر می شود؛ مثلاً بعد از «م» می تواند «ی» بیاید ولی «ل» نمی تواند قرار گیرد. به این ترتیب، اگر پیشوندی در کلمه فعل وجود داشته باشد، جدا و سپس ذخیره می شود. در ضمن، در هر یک از مراحل جداسازی، هرگاه لازم باشد، جست و جویی در بانک بن ها انجام می گیرد.

۴-۳-۳. جدا کردن تک تک حروف از انتهای کلمه و بررسی پسوندها

در این مرحله، اولین حرف سمت چپ - انتهای - کلمه جدا و سپس بررسی می شود که آیا این حرف می تواند جزء آخرین حروف پسوندهای فعل فارسی باشد یا خیر؛ مثلاً حرف «م» می تواند آخرین حرف پسوند (مثل - یم) باشد، ولی مثلاً «ر» آخرین حرف هیچ پسوند فعلی در زبان فارسی نیست. به همین ترتیب، انتخاب های بعدی محدودتر می شود؛ مثلاً قبل از «م» می تواند «ی» بیاید ولی «پ» نمی تواند قرار گیرد. به این ترتیب، اگر پسوندی در کلمه فعل وجود داشته باشد، جدا و سپس ذخیره می شود. در هر یک از مراحل جداسازی، در مواقع لزوم، جست و جویی در بانک بن ها انجام می گیرد.

۴-۳-۴. عملیات میانی برای تشخیص گروه‌های فعلی

برخی از فعل‌ها در فارسی از بیش از یک کلمه تشکیل می‌شود، مانند «رفته است»، «گفته خواهد شد»، «داشتید می‌رفتید»، «فرو خورده شد»، ...

در این گونه افعال - که تعداد زیادی از حالت‌های فعل را دربر می‌گیرد - تمامی حالات بررسی و برای رایانه تعریف شده است؛ به این معنی که با استفاده از مشخصات صوری، همه حالات، حروف، فعل‌های کمکی، پیشوندها و پسوندها، در قبل و بعد از افعال مشخص و کدنویسی شده است، در نتیجه همه افعال چندکلمه‌ای مثل موارد بالا و یا افعال پیشوندی تشخیص داده می‌شود. حتی در مواردی که بین کلمات فعل جدایی می‌افتد، مانند «داشتیم به تهران می‌رفتیم»، با تعریف متغیرهای گوناگون که عمر معینی دارد (نگاه کنید به: کد برنامه)، این موارد برای برنامه قابل تشخیص شده است.

نموداری که در پایان مقاله آمده، تنها نمونه بسیار کوچک و ساده‌ای است از موارد مختلفی که در فعل فارسی وجود دارد، آن هم فقط در افعال ساده. بدیهی است که برای هر یک از افعال گروهی، چندین نمودار از این دست لازم است. چنان‌که ملاحظه می‌شود، در این نمودار فقط حالات مختلف افعالی که به «د» ختم می‌شود و پسوندهای مختلف و ساختمان آنها - بدون پیشوندها - نمایش داده شده است. برای نمونه، اگر اولین مورد را دنبال کنیم، مثلاً برای فعل «خوردن»، خواهیم داشت: خور + که می‌شود خورد، و الی آخر.

برای گروه‌های فعلی به‌ازای هر فعلی که در ساختمان این افعال به کار می‌رود مانند «شدن» برای مجهول و یا «خواستن» برای فعل‌های آینده و یا اسم مفعول در فعل‌های مجهول یا ماضی نقلی و بعید و...، یک مرحله جداگانه نوشته شده است. برای مثال، در مورد فعل «داشت خورده می‌شد»، برنامه ابتدا فعل «داشت» را تشخیص می‌دهد و اعلام می‌کند. از آنجا که فعل «داشتن» از جمله فعل‌های شرکت‌کننده در گروه‌های فعلی است، متغیری فعال می‌شود و آن را ذخیره می‌کند؛ سپس کلمه «خورده» را به سبب اینکه فعل نیست، تشخیص نمی‌دهد و از آن می‌گذرد. آنگاه با رسیدن به کلمه «می‌شد»، آن را به عنوان فعل تشخیص می‌دهد و اعلام می‌کند. در اینجا برنامه به مرحله مخصوص فعل «شدن» می‌رود. طبق

نرم افزار تشخیص فعل در زبان فارسی ۴۱

تعریف، برنامه حرف آخر کلمه ماقبل را - که همیشه در این برنامه ذخیره می شود - جدا می کند. اگر این حرف «ه» نباشد، عملیات پایان می پذیرد، متغیرها خاموش می شوند، و برنامه به کلمه بعد می رود؛ ولی اگر «ه» باشد، بقیه آن کلمه در بانک بن های ماضی جست و جو می شود، اگر پیدا نشود، عملیات مثل مرحله قبل تمام می شود، ولی اگر پیدا شود، آنگاه سه کلمه «داشت»، «خورده»، و «می شد» به هم می پیوندند و به صورت یکجا در برنامه اعلام می شود.

۴-۴. اعلام مشخصات در صورت "فعل" بودن کلمه

با مشخص شدن پیشوندها و پسوندها - اگر کلمه فعل باشد و پیشوند یا پسوند هم داشته باشد - و همچنین بن ماضی یا مضارع فعل، متغیرهای گوناگونی فعال می شود. مثلاً با پیداشدن بن ماضی، متغیری به نام Past و یا با پیداشدن پیشوند «می»، متغیری به نام Mi فعال می شود. در این هنگام، قسمتی از برنامه فعال می شود که وظیفه آن اعلام مشخصات فعل است؛ به این ترتیب که ابتدا تمامی پیشوندها، پسوندها، و بن های یافت شده که ذخیره شده اند، به ترتیب و با علامت «+» در یک کادر فهرست^۱ اعلام می شوند، مانند:

«خورد + ه + خواه + د + شد» برای فعل «خورده خواهد شد»؛

سپس بر مبنای متغیرهای فعال شده، مشخصات فعلی یافت شده در کادر فهرست دیگری اعلام می شود. یعنی مثلاً فعال شدن Mi باعث چاپ شدن استمراری، فعال شدن همزمان صفت مفعولی به علاوه فعل «شدن» نماینده مجهول، و فعال شدن عنصر سببی ساز نشانگر فعل سببی می شود، و الی آخر. بنابراین، برای فعل نمونه بالا چنین پیامی چاپ می شود:

فعل آینده ساده مجهول صیغه سوم شخص مفرد از مصدر خوردن

۴-۵. تشخیص و بررسی کلمه بعد

در این قسمت، برنامه به مرحله ۴-۱ بازمی‌گردد و عملیات جدید را از سر می‌گیرد.

۵. پس‌ویرایش

ساختن نرم‌افزارهایی که به صورت کامل و بدون خطا متن‌ها یا پاره‌های زبانی را تحلیل و یا ترجمه کند، از دهه ۴۰ میلادی همواره آرزوی بشر بوده است. اما این آرزو حتی در زمینه متن‌های تخصصی با واژه‌ها و ساختارهای محدود نیز هیچ‌گاه محقق نشده است؛ و همواره وجود درصدی از خطا در برنامه‌های رایانه‌ای زبانی، طبیعی تلقی می‌شود. به همین جهت، نرم‌افزارهای طراحی شده یا به پیش‌ویرایش^۱ و یا به پس‌ویرایش نیاز دارد؛ به این معنی که مواردی که موجب خطا یا غیرقابل پردازش صحیح را کاربر پیش از دادن متن به رایانه به طور خاص - قابل فهم برای رایانه - علامتگذاری یا پس از پردازش ویرایش می‌کند.

در برنامه مورد بحث ما نیز مواردی مانند منضم کردن جزء غیرفعلی افعال مرکب، فعل‌های ناگذر و موارد مشکل‌ساز دیگر به پس‌ویرایش سپرده می‌شود.

۶. ویژگی‌های برنامه

الف - سرعت بالا: به دلیل کوچک بودن پایگاه داده‌ها - که از فقط حدود ۴۵۰ عضو تشکیل شده است - امکان جست‌وجوهای بسیار سریعی فراهم می‌شود که در نهایت به سرعت زیاد برنامه می‌انجامد. طراحی این نرم‌افزار بر مبنای داده‌های زبان‌شناختی و ویژگی‌های ساختاری زبان فارسی نیز سبب نوشتن کمترین کد با پوشش دادن بیشترین حالات ممکن شده است که خود این امر هم سرعت برنامه را بالا می‌برد.

برای روشن شدن مسئله، مصدر ساده «خوردن» را در نظر می‌گیریم. این مصدر در زمان‌های مختلف با احتساب صورت‌های مجهول، سببی، امر، و حالات مثبت و

نرم افزار تشخیص فعل در زبان فارسی ۴۳

منفی در شش صیغه - بدون احتساب اینکه فعل های «گشتن» و «گردیدن» هم می توانند در صورت های مجهول به کار روند، و البته این امر در برنامه گنجانده شده است - بیش از ۵۴۰ حالت مختلف به خود می گیرد. حال با توجه به اینکه مصدرهای رایج ساده فارسی بیش از ۴۵۰ عدد است، تنها برای افعال ساده - غیرمرکب، غیرپیشوندی،... - به پایگاه داده ای شامل بیش از ۲۵۰ هزار عضو نیاز است! در صورتی که این برنامه با یک پایگاه داده ۴۵۰ عضوی قادر به تشخیص تمامی افعال گروهی نیز است؛ و این خود نشانگر برتری نرم افزارهای زبانی است که برمبنای دستاوردهای زبان شناختی طراحی می شوند.

ب - هوشمند بودن: همان طور که در قسمت الف توضیح داده شد، این برنامه برمبنای جست و جو قرار ندارد؛ بلکه با تحلیلی ساختارزی، مرحله به مرحله، به شناسایی هوشمندانه اجزای کلمه اقدام و با جداسازی و تجزیه این اجزا، نوع کلمه را مشخص می کند. حتی در بخش اعلام مشخصات فعل، این مشخصات برای افعال گوناگون از پیش نوشته نشده است؛ بلکه با فعال شدن متغیرهای گوناگون، توضیحات بخش اعلام مشخصات فعل مرحله به مرحله، کامل می شود، که این امر را می توان به عنوان یک نمونه از شبیه سازی عملکرد مغز برای تشخیص مقولات زبانی در نظر گرفت.

۷. مشکلات کار

این بخش را نگارنده از آن جهت در مقاله گنجانده که در صورتی که پژوهشگران دیگری بخواهند نقایص برنامه موجود را رفع و آن را بهینه کنند و یا نرم افزارهای مشابهی برای دیگر مقولات زبانی طراحی کنند، پیش تر از برخی مشکلات کار آگاه باشند.

الف - هزینه زیاد: چنین پروژه هایی در حد خود بسیار پرهزینه تر از تحقیق در زمینه های دیگر است. هزینه های تهیه رایانه و نرم افزارهای گوناگون، و آموزش را می توان به عنوان هزینه های اصلی این گونه تحقیقات نام برد.

ب - وقت گیر بودن: چنین تحقیقاتی بسیار وقت گیر است و به شکیبایی فراوان در برابر ناکامی های متعدد، جواب نگرفتن از یک شیوه و عوض کردن آن در میانه راه،

و... نیاز دارد.

پ - کار گروهی: در این‌گونه تحقیقات، منطقی‌ترین شیوه، کارکردن گروهی زبان‌شناس مطلع از دانش رایانه‌ای و افراد متخصص رایانه و آشنا با مسائل زبان‌شناسی است؛ که یافتن، ترغیب‌کردن، گردآوردن، و به‌پایان‌رساندن کار این جمع، بس دشوار - و در عین حال بسیار ارزشمند - است.

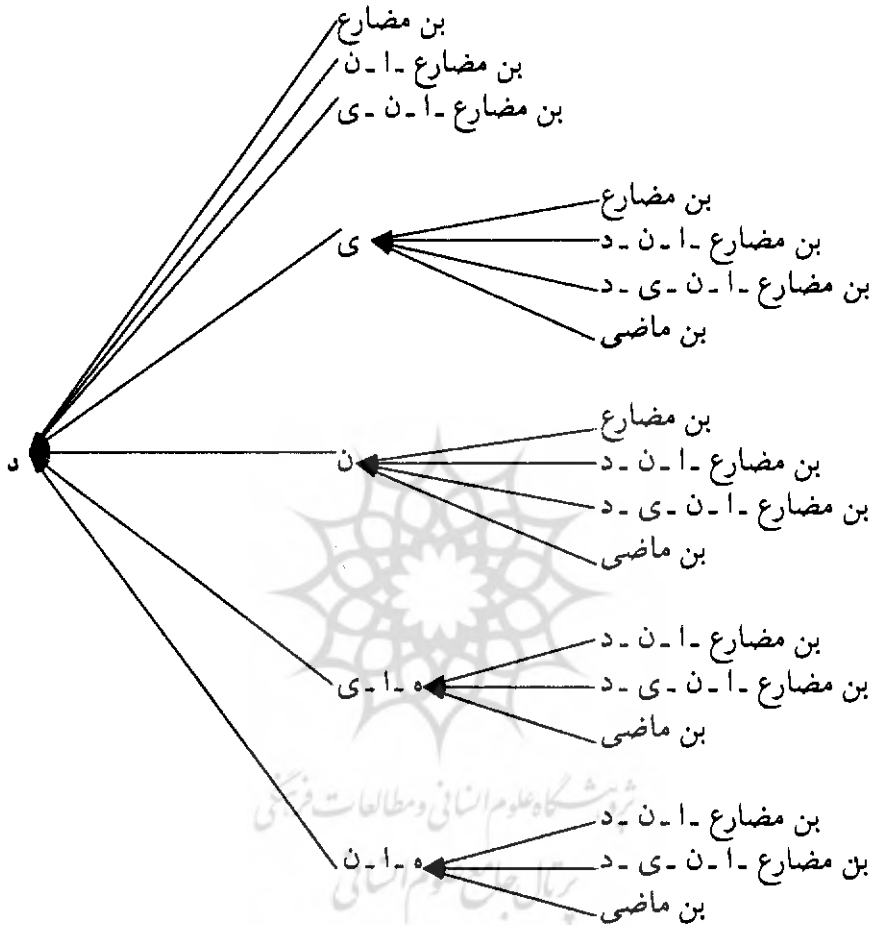
۸. نتیجه‌گیری

چنان‌که ملاحظه شد، ساخت نرم‌افزارهای تحلیل خودکار زبان، کاری دشوار ولی ممکن است. با پیشرفت روزافزون رایانه‌ها در ابعاد مختلف از نظر سرعت، ظرفیت، زبان‌های برنامه‌نویسی، و مانند آن، و همچنین دستاوردهای هرچه بیشتر زبان‌شناسی، و نیز سرمایه‌گذاری‌های مالی و انسانی کافی در این بخش، می‌توان امیدوار بود که در آینده نزدیک، شاهد تحقق رؤیای شصت‌ساله بشر در ساخت ماشین ترجمه باشیم.

مورد دیگری که اهمیت آن کمتر از مطلب فوق نیست، تأیید پیچیدگی و گستردگی قوای زبانی بشر است؛ زیرا همان‌گونه که شرح آن رفت، تنها برای تشخیص مقوله فعل در زبان فارسی گونه نوشتاری معیار، به حدود ۴۰۰۰ خط برنامه نیاز است؛ که خود متکی به میلیون‌ها خط برنامه‌های جانبی از پیش طراحی شده است. اکنون می‌توان تصور کرد که اگر بخواهیم سیستمی طراحی کنیم که همه مقولات دستوری یک زبان و شاید مقولات مشترک و جهانی زبان‌های مختلف را بررسی کند، به چه حجمی از داده‌ها و برنامه‌ها نیاز داریم.

نرم افزار تشخیص فعل در زبان فارسی ۴۵

نمودار پیوندهای فعل در زبان فارسی که به حرف «ن» ختم می‌شوند



کتابنامه

- باطنی، محمدرضا. ۱۳۷۰. توصیف ساختمان دستوری زبان فارسی. چاپ چهارم. تهران: انتشارات امیرکبیر.
- دانش کار آراسته، پویان. ۱۳۸۱. نرم افزار تشخیص فعل در زبان فارسی. پایان نامه کارشناسی ارشد. تهران: دانشگاه علامه طباطبائی.
- دبیرمقدم، محمد. ۱۳۶۷. «ساخت‌های سببی در زبان فارسی»، زبان‌شناسی، س ۵، ش ۱، تهران: مرکز نشر دانشگاهی.
- _____ . ۱۳۷۴. «فعل مرکب در زبان فارسی». زبان‌شناسی، س ۱۲، ش ۱ و ۲. (تاریخ انتشار: ۱۳۷۶)، تهران: مرکز نشر دانشگاهی.

۴۶ فرهنگ، ویژه زبان‌شناسی

عاصی، مصطفی. ۱۳۷۱. «نقش ترکیب در گسترش واژگان زبان فارسی با نگرشی بر آثار نظامی گنجوی»، فرهنگ، کتاب دهم.

_____ . ۱۳۷۲. «کاربرد کامپیوتر در زبان‌شناسی و فرهنگ‌نگاری»، پژوهشگران، ۴، تهران: مؤسسه مطالعات و تحقیقات فرهنگی.

_____ . ۱۳۷۳. «پردازش متن فارسی با کامپیوتر»، پژوهشگران، ۱۲، تهران: مؤسسه مطالعات و تحقیقات فرهنگی.

مشکوة‌الدینی، مهدی. ۱۳۷۴. دستور زبان فارسی بر پایه نظریه گشتاری. چاپ سوم. مشهد: دانشگاه فردوسی مشهد.

ناتل خانلری، پرویز. ۱۳۷۰. دستور زبان فارسی. چاپ دوازدهم. تهران: انتشارات توس.

_____ . ۱۳۷۲. دستور تاریخی زبان فارسی. تهران: انتشارات توس.



پژوهشگاه علوم انسانی و مطالعات فرهنگی
پرتال جامع علوم انسانی