



مرکز منطقه‌ای اطلاع‌رسانی علوم و فناوری

بسمعه تعالی
وزارت علوم، تحقیقات و فناوری
مرکز منطقه‌ای اطلاع‌رسانی علوم و فناوری

گزارش عملکرد سالانه
گروه پژوهشی زبان‌شناسی رایانه‌ای
۱۳۸۸

گروه طرح و برنامه - خرداد ۱۳۸۹

گروه پژوهشی زبان شناسی رایانه ای

مقدمه:

رایانه اگرچه از دستاوردهای غرب است اما در دو دهه اخیر آنچنان در جوامع مختلف از جمله جامعه ما رخنه کرده که امروزه در بعضی از مراکز حتی تصور اینکه روزی رایانه در آن مرکز نباشد، بسیار مشکل است. با وجود این هنوز تنها تعداد محدودی از زبانها با رایانه تطبیق داده شده و در آن قابل استفاده می باشند و زبان فارسی از این جهت در ابتدای راه است.

زبانشناسی رایانه ای به قسمتی از زبانشناسی گفته می شود که مسائل مربوط به زبان و رایانه را بطور عمومی مورد بحث و بررسی قرار می دهد و نشان می دهد که زبانشناسی چگونه می تواند به حل مسائل موجود در مورد تطبیق خط و زبان در رایانه، کمک نماید.

اعضاء این گروه سعی دارند که با اجرای طرحها و نشر کتابها و مقالات و سخنرانیها و ... بتوانند در برطرف کردن مشکلات موجود در خصوص هم آهنگی خط و زبان فارسی با رایانه کمک کرده و در ایجاد یک نگرش علمی نسبت به آن موثر واقع شوند.

در زیر به تعدادی از موارد کاربرد علم زبانشناسی رایانه ای اشاره می شود:

۱- ساخت ریشه ساز زبان فارسی: ریشه سازی یکی از ابزارهایی است که در بازیابی اطلاعات در برخورد با مسئله عدم انطباق واژگان مورد استفاده قرار می گیرد. (یعنی عدم تطبیق واژه های پرس و جو با واژه های مدرک) به فرایند حذف هر پسوند از واژه ها و تبدیل این واژه ها به ریشه های آنها ریشه سازی گفته می شود. برای مثال ریشه سازی واژه "ورزش" را به شکل "ورز" (ریشه مضارع فعل ورزیدن) به وجود می آورد. ریشه سازها عموماً متناسب با هر زبان خاص تهیه می شوند. طراحی ریشه سازها مستلزم خبرگی زبانشناسی در زبان و درک نیازهای مرتبط با بازیابی اطلاعات می باشد.

۲- ساخت ماشین های ترجمه: ترجمه ماشینی یکی از نخستین استفاده هایی است که از پردازش زبان طبیعی به عمل آمده است. برای ساخت اینگونه ماشینها استفاده از تئوری های نحوی و معنایی موجود در زبانشناسی ضروری است. در این زمینه زبان شناسان می توانند در حیطه های: **ایجاد سیستم های ترجمه تماماً خودکار، ترجمه ماشینی نیازمند به انسان، ترجمه با کمک ماشین، ساخت ویرایش گره های املایی، ایجاد سیستم های ترجمه روی خط، ایجاد اصطلاحنامه های روی خط و سرانجام ماشین های همگرا و واگرا** مفید باشند.

نکته حائز اهمیت آنکه برای ساخت ماشین های ترجمه از نظریه اطلاعات و رمزنگاری نیز استفاده می شود که این موضوع ارتباط حیطه های اطلاع رسانی با زبانشناسی رایانه ای را نشان می دهد. بدیهی است مراکز اطلاع رسانی برخوردار از این مزیت قادر خواهند بود اطلاعات مورد نیاز کاربران را به بیش از یک زبان ارایه نمایند.

۳- بازیابی اطلاعات: اساس کار مراکز و پایگاههای اطلاع رسانی بر جستجو و بازیابی اطلاعات بنا شده است. زیر بخشی از زبانشناسی رایانه ای که به **معناشناسی رایانه ای** موسوم است تاکنون توانسته است کیفیت این گونه سیستم ها را به نحو مطلوبی ارتقاء دهد. از **معناشناسی منطقی** نیز استفاده های فراوانی به عمل می آید.

به اختصار می توان گفت که اساس کار سیستم های بازیابی اطلاعات اسناد موجود نیست، بلکه بازنمودهای معنایی است و ایجاد این بازنمودها کاری است در حیطه

معناشناسی رایانه ای. دو مورد از مشکلات عمده موجود در سیستم های بازیابی اطلاعات به ابهام و تشابه واژگانی مربوط است که رفع این مشکلات نیز در حیطه زبانشناسی رایانه ای است.

منطق توصیفی، شبکه های معنایی و صورت گرای قالب- محور مبنای کار سیستمهای اطلاع رسانی در حیطه معناست.

۴- بازشناسی گفتار و سیستم های گفتار به متن: بازشناسی گفتار بر اساس مدلهای زبانشناختی بنا شده که از آن جمله می توان به مدل مارکو اشاره نمود. دستور حالتیهای محدود چامسکی نیز کاربرد فراوان دارد. در سیستم های بازشناسی گفتار شناخت امواج صوتی، چگونگی تعبیر و تفسیر آنها، رمزگذاری و رمزگشایی و چگونگی تغییر مشخصه های آوایی مسایلی عمده و اساسی می باشند که یک آشناس یا واج شناس از عهده آن بر می آید. (آشناسی و واج شناسی یکی از زیر شاخه های علم زبانشناسی است). تاکنون در زمینه تلفظ و املا الگوهای مختلفی ارائه گردیده است.

۵- برجسب زنی کلمات: این حیطه نیز از حیطه های مطرح در علم اطلاع رسانی است و به یکی دیگر از زیر شاخه های علم زبانشناسی که نحو نام دارد مربوط است. شناخت مقوله های مختلف دستوری و چگونگی ایجاد ارتباط و تمایز بین دسته های مختلفی از کلمات به اطلاعات تخصصی در حیطه نحو نیازمند است.

۶- سیستم های متن به گفتار: حیطه رایج دیگری در علم اطلاع رسانی است که از واج شناسی و آشناسی رایانه ای استفاده های زیادی به عمل می آورد. شناخت اندام های صوتی، آوانویسی آوای زبانی، استفاده از الگوها و قواعد واجی و واج شناسی در سیستم های متن به گفتار همگی نیازمند به فردی است که در حیطه واج شناسی به مطالعه و تحقیق پرداخته باشد.

در ادامه برای رعایت اختصار کاربردهای عمده دیگر این حیطه فهرست وار ارائه می گردد.

۷- واژه سازی: جستجو و پیشنهاد معادل های مناسب برای واژه های فاقد معادل فارسی (البته برای این کار لازم است با هماهنگی و یا همکاری فرهنگستان زبان و ادب فارسی اقدام گردد).

۸- نمودارهای N: از این نمودارها برای شمارش اعداد در متن و ... استفاده می شود و مطالب موجود تا حد زیادی به ریاضیات و علم رایانه نیز مربوط است.

۹- دستورهای بافت وابسته و مستقل از بافت: سیستم های بازیابی با استفاده از این گونه الگوها می توانند به جستجو یا تجزیه متن دست بزنند که از آن جمله می توان به تجزیه های از کل به جز و از جز به کل اشاره نمود.

۱۰- تجزیه و تحلیل معنایی

۱۱- گفتمان: که در آن در خصوص چگونگی تعبیر و تفسیر ضماین و یافتن مرجع آنها و نکات مرتبط دیگر بحث می شود که در ماشینهای ترجمه و ... کاربرد فراوان دارد.

۱۲- تبدیل متن به گفتار

۱۳- تبدیل گفتار به متن

۱۴- خلاصه سازی خودکار

۱۵- موتورهای کاوش هوشمند

اهداف مهم گروه پژوهشی زبانشناسی رایانه ای بشرح زیر می باشد:

۱- مطالعه ساختار و واژگان های علمی و پژوهشی جهت دریافت توصیفگر های مناسب

۲- مطالعه و بررسی علوم اطلاع رسانی و کتابداری و زبانشناسی و افزایش بهره وری آنها در سطح کشور و منطقه

۳- مکانیزه کردن نمایه سازی و چکیده نویسی مدارک علمی کشور و منطقه

۴- مطالعه در ساختار زبان فارسی و ارتباط منطقی واژه ها و پراکنش آماری آنها

۵- مطالعه در ساختار ترجمه در زبانهای مختلف با استفاده از رایانه و منطق زبانشناسی

۶- افزایش بهره وری اطلاع رسانی با استفاده از علم زبانشناسی در سطح کشور و منطقه

۷- تعیین راهبردهای جستجو با استفاده از اصلاحنامه ها و تزاروسها و واژه نامه های کامپیوتری

۸- بررسی چگونگی تثبیت اصطلاحات فارسی و کاربرد صحیح آن زبان بعنوان یک زبان علمی

۹- بررسی و رفع مشکلات جستجو و بازیابی اطلاعات رایانه ای با استفاده از منطق و علم زبانشناسی

۱۰- بررسی ارتباط زبانشناسی و بازیابی اطلاعات با الگوهای منطقی زبانشناسی

۱۱- بررسی املائی هوشمند از طریق ارتباط علم زبانشناسی و رایانه

- ۱۲- مطالعه ساختار و صرف و نحو زبانهای طبیعی با استفاده از علم منطق
زبانشناسی
- ۱۳- بررسی چگونگی ساختار (میانجی) واسطه های زبانهای طبیعی برای سیستمهای
رایانه ای
- ۱۴- بررسی راهبردهای تبدیل گفتار به نوشتار و بالعکس با استفاده از علم و منطق
زبانشناسی
- ۱۵- مطالعه خلاصه سازی خودکار در رایانه با استفاده از علم و منطق زبانشناسی
- ۱۶- بررسی چگونگی رو در رویی (ارتباط) زبان طبیعی و سیستمهای رایانه ای
- ۱۷- برجسب زنی کلمات جهت بازشناسی رایانه ای، با استفاده از منطق زبانشناسی
- ۱۸- بررسی ساختار دستورهای بافت وابسته و مستقل از بافت جهت تعیین الگوهای
برای بازیابی اطلاعات.

طرحهای تحقیقاتی
پایان یافته در سال ۱۳۸۸

مفردسازی ماشینی صیغه های جمع اسامی
در زبان فارسی

مجری: شاپوررضا برنجیان

تاریخ شروع: ۸۶/۶/۱۴

تاریخ پایان: ۸۸/۸/۳۰

چکیده:

اهمیت بررسی و تحقیقات در حیطه زبان فارسی بویژه در برخورد این زبان و خط آن با علم رایانه در ماشینهای ترجمه، ریشه یاب و ... بر کسی پوشیده نیست، اگر چه تاکنون پژوهشهایی نیز در این زمینه صورت پذیرفته، اما هنوز موارد بسیاری یافت می شود که هیچ گونه بررسی درباره آنها انجام پذیرفته و تحقیقات در این زمینه ضروری بنظر می رسد. مفردسازی صیغه های جمع اسامی نیز یکی از این موارد می باشد.

مسئله اسم فارسی و ویژگی های آن در برابر موتورهای جستجو در شبکه های کامپیوتری نیز سالهاست که مورد توجه اطلاع رسانی قرار گرفته و این امر کم و بیش در کتابخانه ها و مراکز اطلاع رسانی مورد بررسی قرار گرفته است، طرح حاضر تلاشی است در جهت هموارتر کردن این راه ناهموار.

طرحهاي در دست اجرا

در سال ۱۳۸۸

ریشه ياب ماضي و مضارع از مصدر
افعال ناگذر در زبان فارسي

مجري: شاپوررضا برنجيان

تاريخ شروع: ۸۷/۱۰/۱

چکیده:

تقریباً تمامی زبانهای موجود در دنیا در هر برهه از زمان بدلیل ماهیت خود زبان اقدام به واژه سازي مي کنند، این عمل از دو طریق (اشتقاق و تصریف) صورت مي گیرد، که اکثر آنها از بن های مضارع یا ماضي افعال ساخته مي شوند. مانند: خواستار، گفتار، کردار، دانش ها و ... و همچنین صرف افعال فارسي (در تمام زمانها) نیاز به مشخص نمودن بن ماضي و مضارع افعال دارند، از این رو ساخت ریشه یا بن های ماضي و مضارع در افعال فارسي در ذخیره سازي و بازیابی اطلاعات، سیستمهای ماشينهاي ترجمه، تبدیل متن به گفتار، گفتار با متن و ... از اهمیت خاصی برخوردار است.

چالش های پردازش زبان طبیعی فارسی

مجری طرح: دکتر حمید علیزاده

زمان شروع: ۸۸/۴/۴

چکیده:

امروزه اهمیت کاربرد اصول و تکنیک های پردازش زبان طبیعی در بازیابی اطلاعات به وضوح اثبات شده است. اما علیرغم آنکه در زبان های پرکاربردی همچون انگلیسی، فرانسه، چینی و ... پژوهش های بسیاری در زمینه چالش های موجود برای بکارگیری اصول پردازش زبان طبیعی انجام شده است. در زبان فارسی جدا از چند کار پراکنده که هرکدام وجهی از کاربردهای این حوزه را در نظر گرفته اند. پژوهشی جامع که ظرفیت ها، امکانات و کاستی های زبان فارسی را برای کاربرد اصول پردازش زبان طبیعی بررسی نماید انجام نشده است. در جهت طراحی و ساخت ابزار پردازش زبان طبیعی نظیر ریشه سازها و ابزار تحلیل انواع کلام، ابتدا باید ساختارهای مربوط به زبان فارسی تحلیل گردد تا چنین ابزارهایی با درک صحیح از نیازمندی ها و امکانات موجود طراحی گردد. طرح پژوهشی حاضر، با بررسی چالش های موجود در پردازش زبان طبیعی فارسی، به روشن سازی مسائل اصلی از قبیل امکانات و ساختار خط فارسی، الگوهای زبانی، وضعیت زبان فارسی در پشتیبانی از نظام هایی مثل مترجم های ماشینی و استخراج اطلاعات پرداخته و گزارش نهایی آن که به شکل یک کتاب ارائه خواهد شد، به عنوان مبنای دانش برای پژوهش های مربوط به این حوزه در گروه زبانشناسی رایانه ای مرکز منطقه ای اطلاع رسانی علوم و فناوری و سایر پژوهشگران استفاده خواهد شد.

تکواژگونه و واجگونه

نویسنده: شاپوررضا برنجیان

ارایه شده در فصلنامه علمی - پژوهشی اباختر پائیز ۱۳۸۸ سال پنجم شماره پیاپی ۱۷ و

۱۸

چکیده:

زبان دستگاهی است که از اصوات تشکیل شده است و کارکرد اجتماعی دارد و زبان شناس ابتدا آنرا به صورت صوری مورد تجزیه و تحلیل قرار داده و تقطیع می کند تا بتواند آن را به عنوان نظامی از نشانه ها تعریف کند. در این مقاله بر آن هستیم تا ضمن اشاره به واحدهای گفتار در زبان فارسی به بررسی تکواژ و واج پرداخته و سپس تعریفی از تکواژگونه و واجگونه ارائه نمائیم.

بررسی کارآمدی روش های موجود در بازیابی اطلاعات بین زبانی فارسی - انگلیسی

نویسنده: دکتر حمید علیزاده

ارایه شده در فصلنامه علوم و فناوری اطلاعات دوره ۲۵، شماره ۱، ۱۳۸۸

مروری بر مدل های نوین دسترسی آزاد به نتایج تحقیقات علمی

نویسنده: دکتر حمید علیزاده

ارایه شده در فصلنامه کتاب ماه کلیات، سال دوازدهم، شماره دوازدهم، ۱۳۸۸

سخنراني ها

محورهاي زبان

سخنران: شاپوررضا برنجيان

تاريخ ارائه: ۸۸/۷/۱۵

مکان ارائه: مرکز منطقه اي

چکیده:

موضوع اصلي چگونگي بازيابي اطلاعات از كامپيوتر است. در بازيابي متني ابتدا كاربر متني را در نظر دارد. مانند: "فرزندان ايران در نيروهاي مسلح خدمت مي‌كنند" كه در زمان بازيابي مي بايست به جاي هر يك از واژه ها در زنجيره ي افقي بتوان واژه هاي هم معني را نيز جايگزين كرد، مانند: "جوان" يا "كودك" بجاي "فرزند" تا بتوان اسناد و مدارك مرتبط با آن موضوع بازيابي گردد. صرف نظر از اصطلاحات، هر يك از اجزاي آن عبارت در زنجيره به نوبت در معرض بازيابي قرار مي گيرند.

بررسي كارآمدي روشهاي موجود در بازيابي اطلاعات بين زباني

سخنران: دكتور حميد عليزاده

زمان سخنراني: ۸۸/۴/۲۳

مکان: مرکز منطقه اي

کاربرد پردازش زبان طبیعی در بازیابی اطلاعات بین زبانی

سخنران: دکتر حمید علیزاده

زمان سخنرانی: ۸۸/۸/۱۸

مکان: تهران، انجمن کتابداری و اطلاع رسانی

مروری بر اصول و مفاهیم بازیابی بین زبانی و کاربرد آن

سخنران: دکتر حمید علیزاده

زمان: بهمن ماه ۱۳۸۸

مکان: دانشگاه شیراز

کارگاههای آموزشی

۱- کارگاه آموزشی آشنایی با ISC در دانشگاه پیام نور شیراز

مجری: دکتر حمید علیزاده

۲- کارگاه آموزشی طلایه داران علم ایران در دانشگاه اصفهان

مجری: دکتر حمید علیزاده

۳- کارگاه آموزشی طلایه داران علم ایران در دانشگاه آزاد مرودشت

مجری: دکتر حمید علیزاده