



گزارش نقشه راه پردازش خط و زبان فارسی

محمد بحرانی

مهر ۹۵

مقدمه

- پردازش زبان طبیعی (NLP) یکی از نیازهای عصر فناوری جهت استفاده بهینه از منابع اطلاعاتی است.
- به دلیل اهمیت حفظ و نگهداری از زبان و خط فارسی در محیط رایانه‌ای نیاز به فعالیت‌های حوزه پردازش زبان طبیعی بیش از پیش احساس می‌شود.
- به رغم تلاش‌های صورت گرفته بر روی پردازش رایانه‌ای زبان فارسی، هنوز در این حوزه فاصله زیادی نسبت به زبان‌های دیگر (مانند انگلیسی) وجود دارد.

مقدمه

- دسته‌بندی کلی فعالیت‌ها در حوزه پردازش زبان طبیعی
 - پیکره‌ها و منابع زبانی
 - ابزارهای پایه پردازش زبان طبیعی
 - ابزارها و سامانه‌های کاربردی پردازش زبان طبیعی
 - سامانه‌های کاربردی پردازش متن
 - سامانه‌های کاربردی پردازش گفتار
 - سامانه‌های کاربردی پردازش تصاویر متنی
 - سامانه‌های ترکیبی

وضعیت موجود و وضعیت مطلوب فعالیت‌ها

- پیکره‌ها و منابع زبانی

وضعیت مطلوب	وضعیت موجود در زبان فارسی	نوع پیکره یا منبع زبانی
پیکره با برجسب موجودیت‌های نامدار در انواع مختلف با ۱۰ میلیون کلمه	- پیکره واحدهای اسمی شرکت آرمان رایان شریف در حال توسعه به ۴ میلیون کلمه با ۶ نوع موجودیت اسمی - پیکره موسسه نور با ۵۰۰ هزار کلمه با ۳ نوع موجودیت اسمی	پیکره با برجسب موجودیت‌های نامدار
پیکره برجسب‌خورده با عبارات هم‌مرج با حداقل ۳۰۰ هزار جمله از متن‌های مختلف	- پیکره PCAC-2008 با حدود ۱۰۰ هزار کلمه (۳۱ متن) دارای برجسب مرجع برای حدود ۲۰۰۰ ضمیر - پیکره با حدود ۱۵۰۰۰ جمله در پژوهشگاه خواجه‌نصیر در حال تهیه است.	پیکره برجسب‌خورده با عبارات هم‌مرج
پیکره برجسب‌خورده با انواع نقش‌های معنایی از متون مختلف با حداقل ۱۰ میلیون کلمه	- دادگان موسسه نور با ۳۰۰۰۰ جمله و حدود ۵۰۰ هزار کلمه	پیکره برجسب‌خورده با نقش‌های معنایی
پیکره محاوره‌ای با ۲ میلیارد کلمه از ژانرها و منابع مختلف	- بخشی از پیکره دکتر بیجن‌خان (حدود ۷ میلیون کلمه) - بخشی از پیکره irBlogs (تقریباً ۵۰۰ میلیون کلمه از وبلاگ‌های فارسی)	پیکره متنی محاوره‌ای
درخت‌بانک نحوی با ۵ میلیون کلمه و ۱۵۰ هزار جمله با موضوعات مختلف	- درخت‌بانک نحوی شریف (۳۰۰۰۰ جمله، ۵۰۰ هزار کلمه) - پیکره PerTreebank (۱۰۲۸ جمله) - در دانشگاه تهران (با ۵۰۰ هزار کلمه) و پژوهشگاه خواجه‌نصیر (با ۱ میلیون کلمه) در حال تهیه است.	پیکره تجزیه‌شده نحوی (treebank)

وضعیت موجود و وضعیت مطلوب فعالیت‌ها

• ابزارهای پایه پردازش زبان طبیعی

وضعیت مطلوب	وضعیت موجود در زبان فارسی	نام ابزار یا سامانه
جعبه ابزار با دقت ۹۸٪ برای هر یک از ابزارها	<ul style="list-style-type: none"> - جعبه ابزار مرکز تحقیقات مخابرات - جعبه ابزار آزمایشگاه پردازش زبان طبیعی دانشگاه شهید بهشتی - جعبه ابزار شرکت آرمان رایان شریف - جعبه ابزار گروه سیحه 	<ul style="list-style-type: none"> - واحدساز و یکسان‌ساز - ریشه‌یاب و لم‌یاب - برچسب‌زن مقوله نحوی
تجزیه‌گر نحوی سازهای با دقت ۸۵٪ تجزیه‌گر نحوی وابستگی با دقت ۹۵٪	<ul style="list-style-type: none"> - تجزیه‌گر وابستگی موسسه نور - تجزیه‌گر نحوی آزمایشگاه پردازش زبان طبیعی دانشگاه شهید بهشتی (تجزیه‌گر سازهای) 	تجزیه‌گر نحوی
سامانه برچسب‌زن نقش معنایی با دقت ۹۵٪ برروی جملات حوزه خبری و با دقت ۸۵٪ برروی جملات تخصصی	کارهای پراکنده دانشگاهی با کیفیت متوسط و قابلیت کار بر روی جملات ساده	برچسب‌زن نقوش معنایی
تجزیه‌گر معنایی با قابلیت کار بر روی متون حوزه‌های موردنظر و با دقت ۸۰٪	کارهای پراکنده و اندک دانشگاهی	تجزیه‌گر معنایی

وضعیت موجود و وضعیت مطلوب فعالیت‌ها

• سامانه‌های کاربردی پردازش متن

وضعیت مطلوب	وضعیت موجود در زبان فارسی	نام ابزار یا سامانه
سامانه مترجم ماشینی با قابلیت ترجمه زبان‌های زنده دنیا به فارسی و بالعکس با معیار بلو معادل ۵۰ (با حداقل معادل بلوی گوگل) برای متون عمومی و تخصصی	<ul style="list-style-type: none"> - سامانه مترجم ماشینی ترگمان با معیار بلو معادل ۳۰ در ترجمه انگلیسی به فارسی و ۳۴ در ترجمه فارسی به انگلیسی برای متون خبری - سامانه مترجم ماشینی فرازین با معیار بلو معادل ۲۴ در ترجمه انگلیسی به فارسی برای متون خبری - سامانه مترجم ماشینی پارس با قابلیت ترجمه انگلیسی به فارسی در حوزه عمومی و ۲۸ حوزه تخصصی 	ترجمه ماشینی
سامانه ویرایشگر با قابلیت تشخیص و پیشنهاد برای تصحیح خطاهای املایی، دستوری و معنایی با دقت میانگین ۸۰٪ و قابل استفاده به صورت افزونه در واژه‌پردازهای رایج	<ul style="list-style-type: none"> - سامانه ویرایشگر وفا، قابلیت تشخیص و پیشنهاد برای تصحیح خطاهای املایی و گرامری و معنایی - سامانه ویراستیار، قابلیت تشخیص و پیشنهاد برای تصحیح خطاهای املایی 	تشخیص و تصحیح خطاهای املایی و گرامری و معنایی
سامانه پرسش و پاسخ خودکار در حوزه‌های مختلف و مشابه با IBM Watson	<ul style="list-style-type: none"> - سامانه پرسش و پاسخ قرآن‌جوی در حوزه علوم قرآنی - کارهای پراکنده دانشگاهی با کیفیت متوسط و در حوزه‌های خاص 	پرسش و پاسخ خودکار
سامانه درک زبان طبیعی در حوزه‌های مختلف با دقت ۹۰٪ (مشابه با Apple Siri)	کارهای پراکنده دانشگاهی با کیفیت متوسط و در حوزه‌های خاص	فهم زبان طبیعی
سامانه بازیابی اطلاعات در حد موتور جستجوی گوگل	<ul style="list-style-type: none"> - سامانه‌های بازیابی اطلاعات موجود در موتورهای جستجوی بومی پارسی‌جو، یوز و ... 	بازیابی اطلاعات
<ul style="list-style-type: none"> - خلاصه‌ساز استخراجی با دقت ۸۵٪ - خلاصه‌ساز چکیده‌ای با دقت ۷۵٪ 	<ul style="list-style-type: none"> - سامانه خلاصه‌ساز متنی ایجاز (آزمایشگاه فناوری وب دانشگاه فردوسی مشهد) با دقت ۴۵٪ به صورت استخراجی 	خلاصه‌سازی متن

وضعیت موجود و وضعیت مطلوب فعالیت‌ها

• سامانه‌های کاربردی پردازش متن (ادامه)

نام ابزار یا سامانه	وضعیت موجود در زبان فارسی	وضعیت مطلوب
استخراج اطلاعات	کارهای اندک و پراکنده دانشگاهی با کیفیت پایین	سامانه استخراج اطلاعات، قابل کار بر روی متون عمومی و تخصصی با کیفیت بالا و معیار F معادل ۸۵٪
تشخیص موجودیت‌های نامدار	- سامانه NER شرکت آرمان رایان شریف، با قابلیت ۶ موجودیت نامدار با دقت ۸۰٪ بر روی متون خبری - سامانه NER موسسه نور، با قابلیت ۳ موجودیت نامدار با دقت ۸۰٪ بر روی متون خبری	سامانه NER با قابلیت تشخیص انواع موجودیت‌های نامدار پرکاربرد با دقت ۹۵٪ بر روی انواع متون
مرجع‌یابی ضمیمه	کارهای پراکنده دانشگاهی با کیفیت متوسط برای گروهی از ضمایر (عمدتاً ضمایر متصل) و معمولاً به صورت نیمه‌خودکار	مرجع‌یابی خودکار ضمیمه با دقت ۹۰٪ برای انواع ضمایر متصل و منفصل
رفع ابهام معنایی کلمات	کارهای پراکنده دانشگاهی با کیفیت متوسط و تعداد اندک کلمات هدف	سامانه رفع ابهام معنایی با دقت ۹۰٪ با حداقل ۲۰۰ کلمه هدف
تحلیل نظرات کاربران	کارهای پراکنده دانشگاهی با کیفیت متوسط و برای حوزه‌های خاص (مانند نظرات کاربران در مورد کالاها و محصولات)	سامانه تحلیل نظر کاربران با قابلیت کار بر روی انواع متون و نظرات و با دقت ۸۰٪
تشابه‌یابی در متون	- سامانه مشابهت‌یابی موسسه نور - سامانه مشابهت‌یابی جهاد دانشگاهی	سامانه مشابهت‌یابی با دقت ۹۵٪ با قابلیت کار بر روی انواع متون

وضعیت موجود و وضعیت مطلوب فعالیت‌ها

• سامانه‌های کاربردی پردازش گفتار

نام ابزار یا سامانه	وضعیت موجود در زبان فارسی	وضعیت مطلوب
بازشناسی گفتار بیوسه (میکروفونی و تلفنی)	- موتور و سرویس بازشناسی گفتار بیوسه نوسا با واژگان بزرگ ۱۳۰ هزار کلمه‌ای با دقت ۹۵٪ برای گفتار رسمی و میکروفون اختصاصی و ۹۰٪ با سایر میکروفون‌ها در محیط بدون نویز - موتور بازشناسی گفتار بیوسه شنوا با واژگان بزرگ ۶۵ هزار کلمه‌ای با دقت ۹۰٪ برای گفتار رسمی میکروفونی در محیط بدون نویز	- موتور و سرویس بازشناسی گفتار بیوسه با واژگان بزرگ، با دقت ۹۵٪ برای گفتار رسمی و محاوره‌ای با انواع میکروفون‌ها و ۸۵٪ برای گفتار تلفنی و بدون افت کارایی در محیط‌های نویزی
بازشناسی گفتار کلمات مجزا (میکروفونی و تلفنی)	موتور بازشناسی گفتار کلمات مجزای میکروفونی و تلفنی شرکت عصرگوش با دقت ۹۸٪ بر روی واژگان ۴۰۰ کلمه‌ای در محیط بدون نویز	موتور بازشناسی گفتار کلمات مجزای میکروفونی و تلفنی شرکت عصرگوش با دقت ۹۸٪ بر روی واژگان ۱۰۰ هزار کلمه‌ای و بدون افت کارایی در محیط‌های نویزی
تبدیل متن به گفتار	- موتور تبدیل متن به گفتار آریانا با واژگان ۱۰۰ هزار کلمه‌ای با معیار MOS معادل ۴.۵ و معیار DRT معادل ۹۰٪ برای خوشایندی و طبیعی بودن با صدای مختلف زن و مرد. دقت ۷۰٪ در خوانش کسرده اضافه و دقت ۸۰٪ در خوانش هینکارها - موتور تبدیل متن به گفتار رسا محصول شرکت گاتا با معیار MOS معادل ۴.۵ و معیار DRT معادل ۹۰٪ برای خوشایندی و طبیعی بودن با صدای مختلف زن و مرد. دقت ۸۰٪ در خوانش کسرده اضافه و هینکارها	موتور تبدیل متن به گفتار بدون محدودیت واژگان با خوشایندی و طبیعی بودن بالا (MOS معادل ۵ و DRT معادل ۹۵٪) با صداهای مختلف زن و مرد و خوانش کاملاً صحیح کسرده اضافه و هینکارها
بازایی صدا - بازایی موسیقی - بازایی گفتار (جستجو در گفتار)	کارهای پراکنده دانشگاهی با کیفیت پایین در زمینه بازایی موسیقی و بازایی مستندات گفتاری	- سامانه بازایی مستندات گفتاری با معیار F معادل ۸۰٪ به صورت query by keyword - سامانه بازایی موسیقی با معیار F معادل ۹۰٪ به صورت query by humming و query by example
واژه‌یابی در گفتار	- سامانه واژه‌یابی شرکت عصرگوش پرداز با دقت ۸۰٪ بر روی ۱۰ کلمه کلیدی - سامانه جویا مربوط به پژوهشگاه خواجه‌نصیر	سامانه واژه‌یابی در گفتار با دقت ۹۵٪ با حداقل ۵۰ کلمه کلیدی

وضعیت موجود و وضعیت مطلوب فعالیت‌ها

- سامانه‌های کاربردی پردازش تصاویر متنی

وضعیت مطلوب	وضعیت موجود در زبان فارسی	نام ابزار یا سامانه
نویسه‌خوان نوری بومی با دقت ۹۸٪ برای انواع فونت‌ها و تصاویر با رزولوشن پایین	- نویسه‌خوان نوری پرشین‌نگار با دقت ۹۵٪ برای فونت‌های رایج - نویسه‌خوان نوری آراکس با دقت ۹۵٪ برای فونت‌های رایج - نویسه‌خوان نوری گوگل	بازشناسی نویسه‌های نوری
- سیستم تشخیص دست‌نوشته آنلاین با دقت ۹۸٪ - سیستم تشخیص دست‌نوشته آفلاین با دقت ۹۰٪	- کارهای پراکنده دانشگاهی در زمینه تشخیص دست‌نوشته آنلاین و آفلاین با کیفیت متوسط در تست‌های آزمایشگاهی - سیستم تشخیص دست‌نوشته آنلاین - سامسونگ با دقت بالا دارای کاربرد در گوشی‌های هوشمند	تشخیص دست‌نوشته

اولویت‌بندی فعالیت‌ها (ارزیابی کیفی)

- پیکره‌ها و منابع زبانی

محصولات منتج	پیش‌نیاز	کسب و کار	پیچیدگی	وضعیت موجود	پیکره
تجزیه‌گر نحوی	- پیکره خام - داده‌گان وابستگی	پایین	بالا	در حال توسعه در طرح جویشگر	پیکره تجزیه‌شده نحوی (treebank)
بازشناسی گفتار محاوره ای		پایین	بالا	نامطلوب	پیکره متنی محاوره‌ای
برچسب‌زن نقوش معنایی	پیکره خام	پایین	بالا	نیمه مطلوب	پیکره برچسب‌خورده با نقش‌های معنایی
شناسایی موجودیت‌های نامدار	پیکره خام	پایین	بالا	در حال توسعه در طرح جویشگر و سایر طرح‌ها	پیکره برچسب خورده با موجودیت‌های نامدار
مرجع‌یابی ضمیر	پیکره خام	پایین	متوسط	در حال توسعه در طرح جویشگر و سایر طرح‌ها	پیکره برچسب خورده با عبارات هم‌مرجع

اولویت‌بندی فعالیت‌ها (ارزیابی کیفی)

• ابزارهای پایه پردازش زبان طبیعی

وضعیت موجود	پیچیدگی کسب و کار	پیش‌نیاز	محصولات منتج
نیمه مطلوب	پایین		کلیه پردازش‌های متنی
نیمه مطلوب	متوسط	واژگان	اکثر پردازش‌های متنی
مطلوب	پایین	پیکره متنی با برچسب مقوله نحوی	- تجزیه‌گر نحوی - ترجمه ماشینی - تشخیص و تصحیح خطای املائی و گرامری و معنایی - پرسش و پاسخ خودکار - خلاصه سازی متن - استخراج اطلاعات - تشخیص موجودیت‌های نامدار - مرجع یابی ضمیر - ...
در حال توسعه در طرح جویشرگر	پایین	- پیکره تجزیه‌شده نحوی با دستور زایشی - برچسب زن مقوله نحوی	- ترجمه ماشینی - تشخیص و تصحیح خطای املائی و گرامری و معنایی - پرسش و پاسخ خودکار - خلاصه‌سازی متن - تشخیص موجودیت‌های نامدار - مرجع یابی ضمیر - استخراج اطلاعات - تحلیل نظرات کاربران - ...
نامطلوب	متوسط	پیکره برچسب‌خورده با نقش‌های معنایی	- تجزیه‌گر معنایی - فهم زبان طبیعی - استخراج اطلاعات
نامطلوب	پایین	- پیکره تجزیه‌شده معنایی - برچسب‌زن نقوش معنایی	- فهم زبان طبیعی - استخراج اطلاعات

اولویت‌بندی فعالیت‌ها (ارزیابی کیفی)

• ابزارها و سامانه‌های کاربردی پردازش زبان طبیعی

وضعیت موجود	پیچیدگی کسب و کار	پیش‌نیاز	محصولات منتج
نامطلوب/در حال توسعه در طرح جویشرگر	متوسط	پایین	- استخراج اطلاعات موجودیت‌های نامدار
نیمه مطلوب/در حال توسعه در طرح جویشرگر	بالا	پایین	- پیکره برچسب‌خورده با موجودیت‌های نامدار - ابزارهای پایه (واحدساز، یکسان‌ساز، چانکر، لم یاب، برچسب زن مقوله نحوی)
نیمه مطلوب/در حال توسعه در طرح جویشرگر	بالا	پایین	- ترجمه گفتار به گفتار
نیمه مطلوب	بالا	پایین	- پیکره موازی دوزبانه انگلیسی-فارسی - پیکره موازی چندزبانه - ابزارهای پایه
نامطلوب	متوسط	پایین	- ابزارهای پایه - واژگان
نامطلوب	متوسط	پایین	- پیکره متنی تحلیل احساسات - ابزارهای پایه
نیمه مطلوب/در حال توسعه در سایر طرح‌ها	بالا	پایین	- سامانه محاوره متنی/گفتاری - ترجمه گفتار به گفتار - جستجوی گفتاری - بازیابی گفتار

اولویت‌بندی فعالیت‌ها (ارزیابی کیفی)

- ابزارها و سامانه‌های کاربردی پردازش زبان طبیعی (ادامه)

وضعیت موجود	پیچیدگی	کسب و کار	پیش‌نیاز	محصولات منتج
نامطلوب	بالا	بالا	- بازشناسی گفتار پیوسته - بازیابی اطلاعات	بازیابی گفتار (جستجو در گفتار)
نیمه مطلوب	بالا	بالا	دادگان متون چاپی	بازشناسی نویسه‌های نوری تبدیل تصاویر متنی به گفتار
نامطلوب/در حال توسعه چویشگر	بالا	پایین	- پیکره برچسب خورده با عبارات هم‌مرجع - ابزارهای پایه (واحدساز، یکسان‌ساز، جانکر، لم یاب، برچسب زن مقوله نحوی)	مرجع‌یابی ضمیر
نیمه مطلوب	بالا	متوسط/بالا	- ابزارهای پایه - بازیابی اطلاعات	تشابه‌یابی در متون
نامطلوب	بالا	بالا	دادگان متون دست‌نوشته	تشخیص دست‌نوشته
نامطلوب	بالا	بالا	- بازشناسی گفتار پیوسته - فهم زبان طبیعی - تولیدکننده زبان طبیعی - تبدیل متن به گفتار	سامانه محاوره گفتاری/متنی

اولویت‌بندی فعالیت‌ها (ارزیابی کمی)

- کمی‌سازی آیتم‌های جدول

◦ آیتم «وضعیت موجود»

- مطلوب/در حال توسعه: ۱
- نیمه مطلوب: ۲
- نامطلوب: ۳

◦ آیتم «پیچیدگی»

- بالا: ۱
- متوسط: ۲
- پایین: ۳

اولویت‌بندی فعالیت‌ها (ارزیابی کمی)

- کمی‌سازی آیتم‌های جدول
 - آیتم «کسب و کار»
 - پایین: ۱
 - متوسط: ۲
 - بالا: ۳
 - آیتم «پیش‌نیاز»: مجموع امتیازهای آیتم «وضعیت موجود» هر پیش‌نیاز (با علامت منفی)
 - آیتم «محصولات»: تعداد محصولات منتج از فعالیت

اولویت‌بندی فعالیت‌ها (ارزیابی کمی)

- ضرایب وزنی هر آیتم
 - وضعیت موجود: ۲
 - پیچیدگی: ۱
 - کسب و کار: ۳
 - پیش‌نیاز: -۱
 - محصولات: ۱

اولویت‌بندی فعالیت‌ها (نتایج)

- ده فعالیت با امتیاز بالاتر (امتیاز ۱۲ به بالا)
 ◦ *RFP* برای موارد قرمز رنگ در حال تهیه است.

۳۰	واحدساز/ یکسان‌ساز
۲۳	ریشه‌یاب
۱۵	برچسب‌زن مقوله نحوی
۱۴	تجزیه‌گر نحوی
۱۴	بازشناسی نویسه‌های نوری
۱۴	تشخیص دست‌نوشته
۱۴	بازیابی گفتار
۱۴	واژه‌یابی در گفتار
۱۳	ترجمه گفتار به گفتار
۱۲	برچسب‌زن نقوش معنایی
۱۲	بازشناسی گفتار کلمات مجزا

با سپاس
؟