



## گزارش نقشه راه پردازش خط و زبان فارسی

محمد بحرانی

مردادماه ۹۵

### عناوین مطالب

- مقدمه
- لیست فعالیت‌های حوزه پردازش زبان طبیعی
- وضعیت موجود و وضعیت مطلوب فعالیت‌ها در زبان فارسی
- اولویت‌بندی انجام فعالیت‌ها برای زبان فارسی

## 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

## لیست فعالیت‌های حوزه پردازش زبان طبیعی

- پیکره‌ها و منابع زبانی
  - پیکره‌ها و دادگان‌های متنی
  - دادگان‌های گفتاری
  - دادگان‌های تصاویر متنی
  - واژگان‌ها و هستان‌شناسی‌ها

## لیست فعالیت‌های حوزه پردازش زبان طبیعی

- پیکره‌ها و دادگان‌های متنی
  - پیکره متنی خام
  - پیکره متنی با برچسب مقوله نحوی
  - پیکره متنی محاوره‌ای
  - پیکره موازی محاوره‌ای - رسمی
  - پیکره تجزیه‌شده نحوی (treebank) با دستور سازه‌ای
  - پیکره تجزیه‌شده نحوی (treebank) با دستور وابستگی
  - پیکره برچسب‌خورده با نقش‌های معنایی
  - پیکره برچسب‌خورده با عبارات هم‌مرجع
  - پیکره موازی دو یا چندزبانه
  - پیکره متنی تحلیل احساس
  - پیکره با برچسب موجودیت‌های نامدار

## لیست فعالیت‌های حوزه پردازش زبان طبیعی

- دادگان‌های گفتاری
  - دادگان گفتاری رسمی-میکروفونی
  - دادگان گفتاری- تلفنی
  - دادگان گفتار محاوره‌ای- میکروفونی
  - دادگان گفتار احساسی
- دادگان‌های تصاویر متنی
  - دادگان متون چاپی برای OCR
  - دادگان متون دست‌نوشته
- هستان‌شناسی‌ها و فرهنگ‌ها
  - وردنت
  - گراف دانش
  - فرهنگ ظرفیت نحوی و معنایی افعال

## لیست فعالیت‌های حوزه پردازش زبان طبیعی

- ابزارهای پایه پردازش زبان طبیعی
  - واحدساز (Tokenizer)
  - یکسان‌ساز (Normalizer)
  - ریشه‌یاب (Stemmer)
  - لم‌یاب (Lemmatizer)
  - تحلیل‌گر صرفی (Morphological Analyzer)
  - برچسب‌زن مقوله نحوی (POS Tagger)
  - تجزیه‌گر نحوی (Syntactic Parser)
  - برچسب‌زن نقوش معنایی (Semantic Role Labeler)
  - تجزیه‌گر معنایی (Semantic Parser)

## لیست فعالیت‌های حوزه 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

## لیست فعالیت‌های حوزه پردازش زبان طبیعی

- سامانه‌های کاربردی پردازش متن
  - ترجمه ماشینی
  - تشخیص و تصحیح خطای املائی و گرامری
  - پرسش و پاسخ خودکار
  - فهم زبان طبیعی
  - بازیابی اطلاعات
  - خلاصه‌سازی متن
  - استخراج اطلاعات
  - تشخیص موجودیت‌های نامدار
  - مرجع‌یابی ضمیر
  - رفع ابهام معنایی کلمات

## لیست فعالیت‌های حوزه پردازش زبان طبیعی

- سامانه‌های کاربردی پردازش متن (ادامه)
  - تحلیل نظرات کاربران
  - تشابه‌یابی در متون
  - تشخیص زبان نوشتار
  - قطعه‌بند موضوعی متن
  - ساده‌ساز متن
  - بسط و پیشنهاد عبارت جستجو
  - تولید زبان طبیعی
  - دسته‌بندی و تشخیص موضوع مستندات متنی
  - استخراج کلمات کلیدی از متن
  - تولید خودکار علائم نقطه‌گذاری در متن

## لیست فعالیت‌های حوزه پردازش زبان طبیعی

- سامانه‌های کاربردی پردازش گفتار
  - بازشناسی گفتار پیوسته (میکروفونی و تلفنی)
  - بازشناسی گفتار کلمات مجزا (میکروفونی و تلفنی)
  - تبدیل متن به گفتار
  - بازیابی صدا
  - واژه‌یابی در گفتار

## لیست فعالیت‌های حوزه پردازش زبان طبیعی

### • سامانه‌های کاربردی پردازش تصاویر متنی

- بازشناسی نویسه‌های نوری (OCR)
- تشخیص دست‌نوشته (آن‌لاین/آفلاین)

### • سامانه‌های ترکیبی

- ترجمه گفتار به گفتار
- سامانه محاوره گفتاری/متنی
- تبدیل تصاویر متنی به گفتار

## وضعیت موجود و وضعیت مطلوب فعالیت‌ها

### • پیکره‌ها و منابع زبانی : پیکره‌ها و دادگان‌های متنی

وضعیت مطلوب	وضعیت موجود در زبان فارسی	نوع پیکره یا منبع زبانی
پیکره متنی خام با ۱۰ میلیارد کلمه از ژانرها و منابع مختلف	- پیکره متنی زبان فارسی (پیکره دکتر بیجن‌خان) با ۱۰۰ میلیون کلمه از منابع مختلف - پیکره همشهری با ۱۷۰ میلیون کلمه از متون خبری - پیکره ایسنا با ۵۵۰ میلیون کلمه از متون خبری	پیکره متنی خام
پیکره برجسب‌خورده با ۱۰۰ میلیون کلمه از ژانرها و منابع مختلف	- پیکره متنی زبان فارسی (پیکره دکتر بیجن‌خان) با ۱۰۰ میلیون کلمه برجسب‌خورده	پیکره متنی با برجسب مقوله نحوی
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
پیکره موازی محاوره‌ای-رسمی با حداقل ۵۰ میلیون کلمه	- در برنامه آینده موسسه نور می‌باشد.	پیکره موازی محاوره‌ای - رسمی
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 با دستور سازهای
درخت‌بانک نحوی با ۵ میلیون کلمه و ۱۵۰ هزار جمله با موضوعات مختلف	- دادگان PerDT (۳۰۰۰۰ جمله، ۵۰۰ هزار کلمه) - دادگان اوپسالا (۶۰۰۰ جمله، ۱۵۰ هزار کلمه) - دادگان DepPerTreebank (۱۰۲۸ جمله)	پیکره تجزیه‌شده نحوی با دستور وابستگی

## وضعیت موجود و وضعیت مطلوب فعالیت‌ها

- پیکره‌ها و منابع زبانی : پیکره‌ها و دادگان‌های متنی (ادامه)

وضعیت مطلوب	وضعیت موجود در زبان فارسی	نوع پیکره یا منبع زبانی
پیکره برچسب‌خورده با انواع نقش‌های معنایی از متون مختلف با حداقل ۱۰ میلیون کلمه	- دادگان موسسه نور با ۳۰۰۰۰ جمله و حدود ۵۰۰ هزار کلمه	پیکره برچسب‌خورده با نقش‌های معنایی
پیکره برچسب‌خورده با عبارات هم‌مرجع با حداقل ۳۰۰ هزار جمله از متن‌های مختلف	- پیکره PCAC-2008 با حدود ۱۰۰ هزار کلمه (۳۱ متن) دارای برچسب مرجع برای حدود ۲۰۰۰ ضمیر - پیکره با حدود ۱۵۰۰۰ جمله در پژوهشگاه خواجه‌نصیر در حال تهیه است.	پیکره برچسب‌خورده با عبارات هم‌مرجع
- پیکره موازی فارسی-انگلیسی با حداقل ۲۰۰ میلیون کلمه از حوزه‌های مختلف - پیکره موازی فارسی و حداقل ۱۰ زبان زنده دنیا با حداقل ۱۰۰ میلیون کلمه	- پیکره دوزبانه فارسی-انگلیسی امیرکبیر (AFEC) با ۲۸ میلیون کلمه - پیکره موازی انگلیسی-فارسی تهران (TEP) با ۶۰۰ هزار جمله از زیرنویس فیلم‌ها - پیکره موازی انگلیسی-فارسی میزان با ۱ میلیون جمله در حوزه ادبیات کلاسیک	پیکره موازی دو یا چندزبانه
پیکره با برچسب‌های حاوی بار احساسی با حداقل ۵۰ هزار جمله در حوزه‌های مختلف	- پیکره سنتی‌پرس (SentiPers) با ۱۱۰۰ جمله برچسب‌خورده، در حوزه نظرات در مورد کالاهای دیجیتال	پیکره متنی تحلیل احساس
پیکره با برچسب 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 برچسب موجودیت‌های نامدار

## وضعیت موجود و وضعیت مطلوب فعالیت‌ها

- پیکره‌ها و منابع زبانی : دادگان‌های گفتاری

وضعیت مطلوب	وضعیت موجود در زبان فارسی	نوع پیکره
دادگان گفتار رسمی میکروفونی با حداقل ۱۰۰۰ ساعت گفتار از حداقل ۵۰۰ گوینده و با میکروفون‌های متنوع در شرایط محیطی مختلف	- دادگان فارس‌دات بزرگ، با حدود ۱۴۰ ساعت گفتار میکروفونی تقطیع‌شده از ۱۰۰ گوینده با برچسب واجی - دادگان فارس‌دات با حدود ۵ ساعت گفتار میکروفونی تقطیع‌شده از ۳۰۰ گوینده با برچسب واجی و آوایی	دادگان گفتاری رسمی-میکروفونی
- دادگان تلفنی محاوره‌ای یک‌طرفه با حجم حداقل ۵۰۰ ساعت با تنوع گوشی تلفن و موبایل و از حدود ۲۰۰ گوینده مختلف - دادگان تلفنی از گفتار رسمی با حجم حداقل ۵۰۰ ساعت با تنوع گوشی تلفن و موبایل و از حدود ۲۰۰ گوینده مختلف	- دادگان فارس‌دات بزرگ تلفنی با ۷۰ ساعت محاوره دوطرفه تلفنی که بخشی از آن تقطیع و برچسب‌دهی شده است. - دادگان فارس‌دات تلفنی با حدود ۵ ساعت گفتار محاوره‌ای یک‌طرفه تلفنی با تقطیع و برچسب‌دهی - دادگان تلفنی شرکت عصرگوشی برداز با حدود ۷ ساعت گفتار رسمی تلفنی	دادگان گفتاری - تلفنی
دادگان گفتار محاوره‌ای میکروفونی با حداقل ۵۰۰ ساعت گفتار از حداقل ۳۰۰ گوینده و با میکروفون‌های متنوع در شرایط محیطی مختلف	کارهای پراکنده و اندک دانشگاهی	دادگان گفتار محاوره‌ای- میکروفونی
دادگان گفتار احساسی با حداقل ۵۰۰۰ جمله در انواع حالات احساسی	- پایگاه داده گفتار احساسی Persian ESD با حدود ۵۰۰ جمله در ۵ حالت احساسی - دادگان گفتار احساسی سهند با ۵۰ جمله در ۵ حالت احساسی	دادگان گفتار احساسی



## وضعیت موجود و وضعیت مطلوب فعالیت‌ها

- پیکره‌ها و منابع زبانی : دادگان‌های تصاویر متنی

وضعیت مطلوب	وضعیت موجود در زبان فارسی	نوع پیکره
دادگان متون چاپی با حداقل ۱۰۰۰۰ صفحه از فونت‌های رایج	دادگان farsioocr با ۲۰۰ صفحه اسکن شده با رزولوشن ۳۰۰ نقطه بر اینچ	دادگان متون چاپی برای OCR
دادگان متون دست‌نوشته با حداقل ۱۰۰۰۰ صفحه دست‌نوشته از حداقل ۱۰۰۰ نویسنده	- بانک اطلاعات حروف گسسته - دست‌نویس فارسی با ۱۰ میلیون تصویر از حروف گسسته - مجموعه ارقام دست‌نویس هدی با حدود ۱۰۰ هزار تصویر دست‌نوشته	دادگان متون دست‌نوشته

## وضعیت موجود و وضعیت مطلوب فعالیت‌ها

- پیکره‌ها و منابع زبانی : هستان‌شناسی‌ها و فرهنگ‌ها

وضعیت مطلوب	وضعیت موجود در زبان فارسی	نوع پیکره
وردنت فارسی در حوزه عمومی با حداقل ۱۰۰ هزار مدخل واژگانی با دقت ۹۵٪ و شامل انواع روابط	- فارس‌نت ورژن ۲ با ۳۰ هزار مدخل واژگانی و با دقت ۹۰٪ - فارس‌نت ورژن ۳ با ۱۰۰ هزار مدخل واژگانی در حال توسعه است. - وردنت فارسی حوزه فاوا با ۳۰ هزار مدخل واژگانی	وردنت
گراف دانش با ۲۰۰ هزار موجودیت نامدار دارای ویژگی‌های هستان‌شناسی	یک گراف دانش اولیه با ۲۰۰ هزار موجودیت نامدار و ۵۰۰ هزار رابطه و بدون دارا بودن هستان‌شناسی در حال توسعه است.	گراف دانش
فرهنگ ظرفیت نحوی و معنایی افعال موسسه نور با حداقل ۱۰۰۰۰ فعل ساده و مرکب به همراه ظرفیت نحوی آنها	فرهنگ ظرفیت نحوی و معنایی افعال موسسه نور با ۴۵۰۰ فعل ساده و مرکب به همراه ظرفیت نحوی آنها	فرهنگ ظرفیت نحوی و معنایی افعال

## وضعیت موجود و وضعیت مطلوب فعالیت‌ها

### • ابزارهای پایه پردازش زبان طبیعی

وضعیت مطلوب	وضعیت موجود در زبان فارسی	نام ابزار یا سامانه
واحدساز با دقت ۹۸٪	<ul style="list-style-type: none"> <li>- جعبه ابزار مرکز تحقیقات مخابرات</li> <li>- جعبه ابزار آزمایشگاه پردازش زبان طبیعی دانشگاه شهید بهشتی</li> <li>- جعبه ابزار شرکت آرمان رایان شریف</li> <li>- جعبه ابزار گروه سجه</li> </ul>	واحدساز
یکسان‌ساز با دقت ۹۸٪	<ul style="list-style-type: none"> <li>- جعبه ابزار مرکز تحقیقات مخابرات</li> <li>- جعبه ابزار آزمایشگاه پردازش زبان طبیعی دانشگاه شهید بهشتی</li> <li>- جعبه ابزار شرکت آرمان رایان شریف</li> <li>- جعبه ابزار گروه سجه</li> </ul>	یکسان‌ساز
<ul style="list-style-type: none"> <li>- ریشه‌یاب با دقت ۹۸٪</li> <li>- لم‌یاب با دقت ۹۸٪</li> <li>- تحلیل‌گر صرفی با دقت ۹۵٪</li> </ul>	<ul style="list-style-type: none"> <li>- جعبه ابزار مرکز تحقیقات مخابرات (ریشه‌یاب)</li> <li>- جعبه ابزار آزمایشگاه پردازش زبان طبیعی دانشگاه شهید بهشتی (ریشه‌یاب، تحلیل‌گر صرفی)</li> <li>- پارس‌مورف، تحلیل‌گر صرفی دانشگاه شریف</li> <li>- جعبه ابزار گروه سجه</li> </ul>	<ul style="list-style-type: none"> <li>- ریشه‌یاب</li> <li>- لم‌یاب</li> <li>- تحلیل‌گر صرفی</li> </ul>
برچسب‌زن مقوله نحوی با دقت ۹۸٪	<ul style="list-style-type: none"> <li>- جعبه ابزار مرکز تحقیقات مخابرات</li> <li>- جعبه ابزار شرکت آرمان رایان شریف</li> <li>- جعبه ابزار آزمایشگاه پردازش زبان طبیعی دانشگاه شهید بهشتی</li> </ul>	برچسب‌زن مقوله نحوی

## وضعیت موجود و وضعیت مطلوب فعالیت‌ها

### • ابزارهای پایه پردازش زبان طبیعی (ادامه)

وضعیت مطلوب	وضعیت موجود در زبان فارسی	نام ابزار یا سامانه
تجزیه‌گر نحوی سازهای با دقت ۸۵٪ تجزیه‌گر نحوی وابستگی با دقت ۹۵٪	<ul style="list-style-type: none"> <li>- تجزیه‌گر وابستگی موسسه نور</li> <li>- تجزیه‌گر نحوی دانشکده برق و کامپیوتر دانشگاه تهران (تجزیه‌گر سازهای)</li> <li>- تجزیه‌گر نحوی آزمایشگاه پردازش زبان طبیعی دانشگاه شهید بهشتی (تجزیه‌گر سازهای)</li> </ul>	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

نام ابزار یا سامانه	وضعیت موجود در زبان فارسی	وضعیت مطلوب
ترجمه ماشینی	- سامانه مترجم ماشینی ترگمان با معیار بلو معادل ۲۰ در ترجمه انگلیسی به فارسی و ۲۴ در ترجمه فارسی به انگلیسی برای متون خبری - سامانه مترجم ماشینی فراژین 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	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
تشخیص و تصحیح خطای املایی و گرامری	- سامانه ویرایشگر وفاء قابلیت تشخیص و پیشنهاد برای تصحیح خطاهای املایی و گرامری - سامانه ویراستیار، قابلیت تشخیص و پیشنهاد برای تصحیح خطاهای املایی	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
پرسش و پاسخ خودکار	- سامانه پرسش و پاسخ قرآن‌جوی در حوزه علوم قرآنی - کارهای پراکنده دانشگاهی با کیفیت متوسط و در حوزه‌های خاص	سامانه پرسش و پاسخ خودکار در حوزه‌های مختلف و با کیفیت بالا
فهم زبان طبیعی	کارهای پراکنده دانشگاهی با کیفیت متوسط و در حوزه‌های خاص	سامانه درک زبان طبیعی در حوزه‌های مختلف و با دقت ۹۰٪
بازایی اطلاعات	- سامانه‌های بازایی اطلاعات موجود در موتورهای جستجوی بومی پارسی‌جو، یوز و ...	سامانه بازایی اطلاعات با معیار F معادل ۸۵٪
خلاصه‌سازی متن	- سامانه خلاصه‌ساز متنی ایجاز (آزمایشگاه فناوری وب دانشگاه فردوسی مشهد) با دقت ۴۵٪ به صورت استخراجی	- خلاصه‌ساز استخراجی با دقت ۸۵٪ - خلاصه‌ساز چکیده‌ای با دقت ۷۵٪

## 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 ضمیر	کارهای پراکنده دانشگاهی با کیفیت متوسط برای گروهی از ضمائر (عمدتاً ضمائر متصل) و معمولاً به صورت نیمه‌خودکار	مرجع‌یاب خودکار ضمیر با دقت ۹۰٪ برای انواع ضمائر 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 مؤسسه 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 متن (ادامه)

نام ابزار یا سامانه	وضعیت موجود در زبان فارسی	وضعیت مطلوب
قطعه‌بند موضوعی متن	کارهای بسیار اندک دانشگاهی بر روی متونی از حوزه‌های خاص	سامانه قطعه‌بند موضوعی متن با قابلیت کار بر روی انواع متون و با دقت ۹۰٪
ساده‌ساز متن	کارهای بسیار اندک دانشگاهی بر روی متونی از حوزه‌های خاص	سامانه ساده‌ساز متن با قابلیت کار بر روی انواع متون و با دقت ۸۵٪
بسط و پیشنهاد عبارت جستجو	سامانه‌های موجود در موتورهای جستجوی بومی	سامانه بسط و پیشنهاد با کیفیت مشابه گوگل
تولید زبان طبیعی	کارهای بسیار اندک دانشگاهی در حوزه‌های خاص	سامانه تولید زبان طبیعی با قابلیت تولید متون طبیعی و روان در زمینه‌های مختلف
دسته‌بندی و تشخیص موضوع مستندات متنی	کارهای پراکنده دانشگاهی با کیفیت متوسط بر روی حدوداً ۱۰ - ۲۰ موضوع	سامانه دسته‌بندی و تشخیص موضوع متون با دقت ۹۰٪ بر روی حداقل ۱۰۰ موضوع مختلف
استخراج کلمات کلیدی از متن	کارهای پراکنده و اندک دانشگاهی با کیفیت متوسط و بر روی حوزه‌های خاص	سامانه استخراج کلمات کلیدی از متون حوزه‌های مختلف با معیار F معادل ۸۵٪
تولید خودکار علائم نقطه‌گذاری در متن	کارهای پراکنده و اندک دانشگاهی با کیفیت متوسط	سامانه تولید خودکار انواع علائم نقطه‌گذاری برای متون حوزه‌های مختلف و با معیار F معادل ۹۵٪

## 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
بازیابی صدا	کارهای پراکنده دانشگاهی با کیفیت پایین در زمینه بازیابی موسیقی و بازیابی مستندات گفتاری	- سامانه بازیابی مستندات گفتاری با معیار F معادل ۸۰٪ به صورت query by keyword - سامانه بازیابی موسیقی با معیار F معادل ۹۰٪ به صورت query by humming و query by example
واژه‌یابی در گفتار	- سامانه واژه‌یاب شرکت عصرگوش پردازش 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

## وضعیت موجود و وضعیت مطلوب فعالیت‌ها

- سامانه‌های ترکیبی

وضعیت مطلوب	وضعیت موجود در زبان فارسی	نام ابزار یا سامانه
سامانه ترجمه گفتار به گفتار با قابلیت ترجمه جملات مختلف زبان‌های گوناگون به فارسی و بالعکس	- مترجم هوشمند همراه شرکت سروش مهر با قابلیت ترجمه ۲۰۰۰ جمله از فارسی به ۹ زبان مختلف و بالعکس - مترجم پارسیا محصول شرکت عصرگوش پرداز با قابلیت ترجمه ۵۰۰ جمله از فارسی به انگلیسی و عربی (هر دو محصول ترجمه را بر اساس translation memory انجام می‌دهند).	ترجمه گفتار به گفتار
سامانه محاوره گفتاری مشابه با Apple SIRI	- سامانه دیالوگ گفتاری نیوشا محصول شرکت عصرگوش پرداز برای کاربردهای تلفن‌گویا و تلفن‌بانک - کارهای پراکنده دانشگاهی با کیفیت متوسط و خاص منظوره	سامانه محاوره گفتاری /متنی
سامانه تبدیل تصاویر متنی به گفتار مشابه با MIT finger reader	- یک سامانه اولیه با قابلیت خواندن یک فونت خاص در دانشگاه تهران در حال توسعه است.	تبدیل تصاویر متنی به گفتار

## اولویت بندی فعالیت‌ها: پیکره‌ها و منابع زبانی

### • اولویت ۱:

- پیکره تجزیه شده نحوی (treebank) با دستور سازه‌ای در حجم ۱ میلیون کلمه
- پیکره متنی محاوره‌ای در حجم ۵۰۰ میلیون کلمه از منابع مختلف
- پیکره موازی محاوره‌ای-رسمی با ۱۰ میلیون کلمه
- دادگان گفتار محاوره‌ای با حجم ۱۰۰ ساعت گفتار از منابع مختلف
- دادگان متون دست‌نوشته با ۱۰۰۰ صفحه
- گراف دانش

## اولویت بندی فعالیت‌ها: پیکره‌ها و منابع زبانی

### • اولویت ۲:

- ادامه فعالیت‌های مذکور در اولویت ۱
- پیکره موازی دوزبانه انگلیسی-فارسی در حجم ۱۰۰ میلیون کلمه
- پیکره موازی چندزبانه در حجم ۴۰ میلیون کلمه و با ۵ زبان زنده دنیا (به غیر از انگلیسی)
- پیکره برچسب خورده با نقش‌های معنایی در حجم ۲ میلیون کلمه
- پیکره با برچسب موجودیت‌های نامدار با ۵ میلیون کلمه
- پیکره برچسب خورده با عبارات هم‌مرجع با ۲.۵ میلیون کلمه
- وردنت با ۱۰۰ هزار مدخل واژگانی
- دادگان متون چاپی برای OCR با ۲۰۰۰ صفحه از ۵ فونت مختلف

## اولویت بندی فعالیت‌ها: پیکره‌ها و منابع زبانی

### • اولویت ۳:

- پیکره تجزیه شده نحوی با دستور وابستگی با ۵ میلیون کلمه
- پیکره متنی خام
- پیکره متنی با برچسب مقوله نحوی
- پیکره تحلیل احساس با ۲۰۰۰۰ جمله

## اولویت بندی فعالیت‌ها: ابزارهای پایه

### • اولویت ۱:

- برچسبزن نقوش معنایی با دقت ۸۰٪
- تجزیه‌گر معنایی برای حوزه‌های خاص با دقت ۶۵٪

### • اولویت ۲:

- ادامه فعالیت‌های اولویت ۱
- تجزیه‌گر نحوی سازه‌ای با دقت ۷۵٪
- ریشه‌یاب با دقت ۹۵٪

### • اولویت ۳:

- ادامه فعالیت‌های اولویت ۲
- واحدساز
- یکسان‌ساز
- برچسبزن مقوله نحوی

## اولویت بندی فعالیت‌ها: سامانه‌های کاربردی

### • اولویت ۱:

- سامانه محاوره گفتاری/متنی برای چند حوزه خاص با دقت ۸۰٪ در فهم زبان طبیعی
- سامانه استخراج اطلاعات و استخراج روابط از متن
- سامانه خودکار مرجع‌یاب ضمیر با دقت ۷۵٪
- سامانه تشخیص موجودیت‌های نامدار با دقت ۸۰٪ برای ۶ نوع موجودیت
- سامانه بازیابی صدا برای مستندات گفتاری با دقت ۶۵٪
- سامانه بازشناسی گفتار پیوسته محاوره‌ای با دقت ۸۵٪
- سامانه تشخیص دست‌نوشته آنلاین/آفلاین با دقت حدود ۸۰٪

## اولویت بندی فعالیت‌ها: سامانه‌های کاربردی

### • اولویت ۲:

- ادامه فعالیت‌های اولویت ۱
- سامانه مترجم ماشینی برای ترجمه زبان‌های مختلف به فارسی و برعکس
- توسعه سامانه پرسش و پاسخ خودکار برای کاربرد در حوزه عمومی
- سامانه خلاصه‌ساز متن به صورت استخراجی و با دقت ۶۵٪
- توسعه سامانه تشخیص و تصحیح خطای املائی و گرامری
- توسعه سامانه‌های تشابه‌یاب متون
- سامانه بازشناسی گفتار کلمات مجزا برای ۵۰۰۰ کلمه مجزا با دقت ۹۵٪
- سامانه واژه‌یابی در گفتار با دقت ۸۰٪ بر روی ۲۵ کلمه کلیدی
- افزایش کیفیت سامانه تبدیل متن به گفتار (TTS) فارسی
- توسعه سامانه بازیابی اطلاعات جویشرهای بومی
- توسعه سامانه بسط و پیشنهاد عبارت جستجو در جویشرهای بومی
- توسعه سامانه‌های OCR برای فونت‌های مختلف



## اولویت‌بندی فعالیت‌ها: سامانه‌های کاربردی

### • اولویت ۳:

- ادامه فعالیت‌های اولویت ۲
- سامانه تحلیل نظرات کاربران در چند حوزه کاربردی
- سامانه تولید خودکار علائم نقطه‌گذاری در متن
- سامانه خلاصه‌ساز متن به صورت چکیده‌ای با دقت ۵۵٪
- سامانه استخراج کلمات کلیدی از متن
- سامانه تولید زبان طبیعی
- سامانه ساده‌ساز متن
- سامانه تشخیص زبان نوشتار
- سامانه قطعه‌بند موضوعی متن
- سامانه دسته‌بندی و تشخیص موضوع مستندات متنی
- سامانه ترجمه گفتار به گفتار با قابلیت ترجمه جملات انگلیسی به فارسی و برعکس
- سامانه تبدیل تصاویر متنی به گفتار

با سپاس

؟