




	<b>عنوان پروژه:</b> فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 شورای عالی اطلاع رسانی
	<b>عنوان زیر پروژه:</b> بررسی ابعاد و تفاوت های پیکره های برجسب داده ای و پیکره های خام در زبان فارسی		
	تاریخ: 1388/04/04	ویرایش: 1/0	

عنوان زیر پروژه:



## **بررسی ابعاد و تفاوت های پیکره های برجسب داده ای و پیکره های خام در زبان فارسی**

	<b>عنوان پروژه:</b> فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 شورای عالی اطلاع رسانی
	<b>عنوان زیر پروژه:</b> بررسی ابعاد و تفاوت های پیکره های برجسب داده ای و پیکره های خام در زبان فارسی		
	تاریخ: 1388/04/04	ویرایش: 1/0	

	<b>عنوان پروژه:</b> فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 شورای عالی اطلاع رسانی
	<b>عنوان زیر پروژه:</b> بررسی ابعاد و تفاوت های پیکره های برجسب داده ای و پیکره های خام در زبان فارسی		
	تاریخ: 1388/04/04	ویرایش: 1/0	

## فهرست مطالب

شماره صفحه	عنوان
4.....	1. مقدمه
8.....	2. پیکره متنی
11.....	3. پیکره های متنی برجسب داده ای و نقش آن در کاربردهای پردازش زبان طبیعی
14.....	1-3. خلاصه سازی
16.....	2-3. ترجمه ماشینی
19.....	4. نتیجه گیری

	<b>عنوان پروژه:</b> فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 شورای عالی اطلاع رسانی
	<b>عنوان زیر پروژه:</b> بررسی ابعاد و تفاوت های پیکره های برجسب داده ای و پیکره های خام در زبان فارسی		
	تاریخ: 1388/04/04	ویرایش: 1/0	

## 1. مقدمه

پردازش زبان طبیعی عبارتست از استفاده از رایانه برای پردازش زبان گفتاری و نوشتاری. پردازش زبان طبیعی می تواند در سطوح مختلف از زبان صورت گیرد. این سطوح را می توان به شکل زیر تقسیم بندی نمود:

1. آواشناسی و صدا شناسی<sup>1</sup> که به تشخیص آواها و صداها و بازشناسی گفتار می پردازد.
2. ریخت شناسی<sup>2</sup> که به ساختارهای کلمات و ریشه یابی واژگان می پردازد.
3. نحو<sup>3</sup> که به ارتباط کلمات به همدیگر و مباحث دستوری آنها در گروه ها و جملات می پردازد.
4. معناشناسی<sup>4</sup> که به ارتباطات معنایی کلمات می پردازد.
5. عمل گرایی<sup>5</sup> که کاربردهای زبان برای رساندن یک مطلب به مخاطب یا مخاطبان، در حالت عملی و یا در نوشتار و گفتار طبیعی می پردازد.
6. مباحثه<sup>6</sup> که به ارتباطات کلی یک زبان فرای یک یا چند جمله خاص می پردازد.

از کاربردهای اصلی پردازش زبان طبیعی می توان موارد زیر را ذکر کرد: (1) خلاصه سازی خودکار<sup>7</sup>، (2) کمک به خواندن زبان های طبیعی دیگر<sup>8</sup>، (3) کمک به نوشتن به زبان های طبیعی دیگر<sup>9</sup>، (4) استخراج

<sup>1</sup> Phonetics and Phonology

<sup>2</sup> Morphology

<sup>3</sup> Syntax

<sup>4</sup> Semantics



<sup>5</sup> Pragmatics

<sup>6</sup> Discourse

<sup>7</sup> Automatic Summarization

<sup>8</sup> Foreign language reading aid

<sup>9</sup> Foreign language writing aid

	<b>عنوان پروژه:</b> فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 شورای عالی اطلاع رسانی
	<b>عنوان زیر پروژه:</b> بررسی ابعاد و تفاوت های پیکره های برجسب داده ای و پیکره های خام در زبان فارسی		
	تاریخ: 1388/04/04	ویرایش: 1/0	

اطلاعات<sup>۱</sup>، ۵) بازیابی اطلاعات<sup>۲</sup>، ۶) ترجمه ماشینی<sup>۳</sup>، ۷) تشخیص واحدهای اسمی<sup>۴</sup>، ۸) تولید زبان طبیعی<sup>۵</sup>، ۹) فهم زبان طبیعی<sup>۶</sup>، ۱۰) نویسه خوان نوری<sup>۷</sup>، ۱۱) تحلیل مرجع دارها<sup>۸</sup>، ۱۲) سیستم سوال، پاسخ<sup>۹</sup>، ۱۳) تشخیص گفتار<sup>۱۰</sup>، ۱۴) مبدل متن به گفتار<sup>۱۱</sup>، ۱۵) نظام های مکالمه گفتاری<sup>۱۲</sup>، ۱۶) ساده سازی متن<sup>۱۳</sup>، ۱۷) تایید متن<sup>۱۴</sup>

با توجه به کاربردهای فوق، برخی مسائل اساسی در پردازش زبان طبیعی را می توان به صورت زیر برشمرد:

1. قطعه بندی گفتار<sup>۱۵</sup>

2. قطعه بندی متن<sup>۱۶</sup>

3. برجسب گذاری نقش کلمه<sup>۱</sup>

Information Extraction<sup>۱</sup>

Information Retrieval<sup>۲</sup>

Machine Translation<sup>۳</sup>

Name Entity Recognition<sup>۴</sup>

Natural language generation<sup>۵</sup>

Natural language Understanding<sup>۶</sup>

Optical Character Recognition - OCR<sup>۷</sup>

Anaphora Reservation<sup>۸</sup>

Question Answering System<sup>۹</sup>

Speech Recognition<sup>۱۰</sup>

Text-to-Speech<sup>۱۱</sup>



Spoken Dialogue System<sup>۱۲</sup>

Text Simplification<sup>۱۳</sup>

Text Proofing<sup>۱۴</sup>

Speech Segmentation<sup>۱۵</sup>

Text Segmentation<sup>۱۶</sup>

	<b>عنوان پروژه:</b> فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 شورای عالی اطلاع رسانی
	<b>عنوان زیر پروژه:</b> بررسی ابعاد و تفاوت های پیکره های برجسب داده ای و پیکره های خام در زبان فارسی		
	تاریخ: 1388/04/04	ویرایش: 1/0	

4. ابهام زدائی از نقش کلمه 2

5. ابهام زدائی نحوی 3

6. هنجار سازی 4

7. تشخیص اعمال گفتاری 5

پردازش زبان طبیعی در حوزه متن یا خط و زبان به پردازش پیکره های متنی<sup>1</sup> به عنوان نمایانگر زبان می پردازد. بنابر این مسائل اصلی در حوزه خط و زبان را می توان به موارد زیر کاهش داد:

1. قطعه بندی متن

2. برجسب گذاری نقش کلمه

3. ابهام زدائی از نقش کلمه

4. ابهام زدائی نحوی

5. هنجار سازی

مسائل فوق تقریباً در تمامی کاربردهای پردازش زبان طبیعی در حوزه خط و زبان مطرح هستند و به عنوان مسائل زیرساختی خط و زبان می باید مورد بررسی قرار گیرند. زبان فارسی به عنوان یکی از زبان های طبیعی نیز از این قاعده مستثنی نیست. از طرفی همان طور که مطرح شد، در حوزه خط و

<sup>1</sup> Part-of-Speech Tagging



<sup>2</sup> Word Sense Disambiguation

<sup>3</sup> Syntactic Disambiguation

<sup>4</sup> Normalization

<sup>5</sup> Speech acts

<sup>6</sup> Text Corpus

	<b>عنوان پروژه:</b> فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 شورای عالی اطلاع رسانی
	<b>عنوان زیر پروژه:</b> بررسی ابعاد و تفاوت های پیکره های برچسب داده ای و پیکره های خام در زبان فارسی		
	تاریخ: 1388/04/04	ویرایش: 1/0	



زبان، پیکره های متنی به عنوان نمادی از زبان هستند که می باید با تحلیل آن ها به استخراج اجزا، قواعد و ساز و کار زبان پی برد.

دادگان مهم و زیرساختی حاصل از تحلیل پیکره های متنی در پردازش زبان طبیعی و حوزه خط و زبان را می توان به موارد زیر تقسیم نمود:

1. پیکره متنی
2. واژگان زبان (واژگان عمومی و تخصصی)
3. گنجوازه یا اصطلاح نامه<sup>1</sup>
4. الگوهای زبان
5. پیکره های برچسب داده ای
6. پیکره های تخصصی

از میان دادگان فوق، تمرکز این مستند بر روی پیکره های متنی و بویژه پیکره های برچسب داده ای است. بنابراین، ابتدا به معرفی پیکره متنی و پیکره های برچسب داده ای پرداخته و سپس نقش پیکره های برچسب داده ای در دو مورد از کاربردهای پردازش زبان طبیعی (خلاصه سازی و ترجمه ماشینی) را برمی شماریم.

<sup>1</sup> Thesaurus

	<b>عنوان پروژه:</b> فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 شورای عالی اطلاع رسانی
	<b>عنوان زیر پروژه:</b> بررسی ابعاد و تفاوت های پیکره های برچسب داده ای و پیکره های خام در زبان فارسی		
	تاریخ: 1388/04/04	ویرایش: 1/0	

## 2. پیکره متنی

در علوم زبان شناسی و زبان شناسی رایانه‌ای، پیکره متنی، حجم بسیاری از متون ساخت یافته آن زبان است. در (Atkins and Clear, 1992) نیز تعریفی که برای پیکره ارائه شده است به این صورت است: حجم زیادی از داده‌های زبانی که براساس معیارهای مشخص برای هدف معینی جمع‌آوری و ذخیره شده‌اند به طوری که نماینده زبان یا گویش مورد مطالعه باشد. به طور کلی در طراحی و تهیه یک پیکره برچسب‌خورده یکی از مهمترین مسائلی که باید مورد توجه قرار گیرد مجموعه برچسب پیکره است که بر اساس هدف غایی پیکره و منظوری که از پیکره مد نظر است حاوی برچسب‌هایی خواهد شد که نیل به آن هدف و منظور را ممکن سازد.

پیکره‌های متنی به منظور تحلیل‌های آماری، صحت‌سنجی فرضیه‌ها و بررسی رخداد یا صحت قواعد زبانی در حوزه‌ای مشخص به کار می‌روند. پیکره‌های متنی به عنوان پایگاه دانش<sup>1</sup> اصلی در زبان‌شناسی رایانه‌ای خصوصاً در حوزه خط و زبان شناخته می‌شوند. پیکره‌های می‌توانند تک زبانی<sup>2</sup> (شامل متون با یک زبان) و یا چند زبانی<sup>3</sup> (شامل متون با چندین زبان) باشند، که مورد اخیر معمولاً برای مقایسه نظیر-به-نظیر ساختار دهی می‌شود و پیکره موازی منطبق<sup>4</sup> نامیده می‌شود.

به منظور مناسب‌سازی پیکره‌های متنی برای کاربرد مفیدتر در مطالعات زبان‌شناسی، پیکره‌ها می‌باید حاشیه‌نویسی<sup>5</sup> شوند. حاشیه‌نویسی عبارت است از تحلیل و افزودن برخی اطلاعات مانند اطلاعاتی در مورد نقش<sup>6</sup> و یا ریشه<sup>7</sup> کلمات موجود در متن به پیکره. پیکره‌های متنی اکثراً در سطح ساختاری

<sup>1</sup> Knowledgebase

<sup>2</sup> Monolingual

<sup>3</sup> Multilingual



<sup>4</sup> Aligned Parallel Corpus

<sup>5</sup> Annotated

<sup>6</sup> Part-of-Speech - POS

<sup>7</sup> Lemma



	<b>عنوان پروژه:</b> فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 شورای عالی اطلاع رسانی
	<b>عنوان زیر پروژه:</b> بررسی ابعاد و تفاوت های پیکره های برجسب داده ای و پیکره های خام در زبان فارسی		
	تاریخ: 1388/04/04	ویرایش: 1/0	

حاشیه نویسی می شوند. پیکره هایی که کاملاً تجزیه<sup>1</sup> و حاشیه گذاری شده باشند، پیکره های تجزیه شده یا بانک درختی<sup>2</sup> گفته می شود. عملاً پیکره های کوچک (بین 1 تا 3 میلیون کلمه) امکان تجزیه کامل را دارند زیرا حصول اطمینان از این که حاشیه نویسی کاملاً صحیح و سازگار است، عملی بسیار پیچیده، زمان گیر و پرهزینه خواهد بود. امکان حاشیه نویسی در سطوح ریخت شناسی، نحوی، معنایی و کاربردی معمولاً برای پیکره های کوچک امکان پذیر است.

از چالش های مطرح در خصوص پیکره های متنی می توان به این نکته اشاره کرد که چه پیکره ای با چه خصوصیتی می تواند به بهترین و بیشترین شکل خصوصیات زبان را بیان سازد. با توجه به تغییر پذیری زبان در زمان، این نکته که چه متونی با چه ویژگی هایی می توانند پیکره ای مناسبی را شکل دهند اهمیت می یابد. این که چه سهمی از پیکره از زبان معیار انتخاب شود، چه میزان آن از زبان فوق معیار باشد و چه بخشی از آن را زبان زیر معیار دربر گیرد از چالش های اساسی در انتخاب و ایجاد پیکره ها است. همچنین تعریف مناسب از زبان معیار، فوق و زیر معیار نیز در این امر تاثیر گذار است و می باید مورد توجه قرار گیرد.



در سه دهه ی گذشته، بیشتر کشورهای صنعتی به ایجاد بانک های زبانی خود پرداخته اند. این کشورها، ابتدا با گردآوری داده ها یا پیکره های متنی، و سپس با سازمان دهی آن ها در پایگاه های داده ها و بانک های اطلاعات زبانی، از شبکه های بین المللی داده های زبانی<sup>3</sup> بهره برداری می کنند.

بسیاری از پژوهش های زبان شناختی و تصمیم گیری ها در برنامه ریزی زبانی، تنها با استفاده از یک پیکره زبانی امکان پذیر است. بنابراین، پیکره متنی پایه ای ترین و یکی از مهم ترین دادگان ورودی در کاربردهای پردازش زبان طبیعی بشمار می رود. در واقع، پیکره های متنی، پایگاه دانش اصلی در زبان شناسی پیکره ای است. تحلیل و پردازش انواع مختلفی از پیکره ها، موضوع مقاله های بسیاری در حوزه ی زبان شناسی رایانشی و پردازش زبان طبیعی بوده است. ایجاد مدل های مارکف از پیکره های متنی، برای اهدافی همچون برجسب زنی نقش کلمه مورد استفاده قرار می گیرد. پیکره های متنی، کاربردهایی خارج از حوزه ی پردازش زبان طبیعی نیز دارند. بعنوان مثال می توان فهرست هایی از فراوانی کلمات را از پیکره ی



<sup>1</sup> Parse

<sup>2</sup> Treebank

<sup>3</sup> International Networks of Linguistic Data

	<b>عنوان پروژه:</b> فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 شورای عالی اطلاع رسانی
	<b>عنوان زیر پروژه:</b> بررسی ابعاد و تفاوت های پیکره های برجسب داده ای و پیکره های خام در زبان فارسی		
	تاریخ: 1388/04/04	ویرایش: 1/0	

متنی استخراج کرده و در آموزش زبان مورد استفاده قرار داد.



	<b>عنوان پروژه:</b> فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 شورای عالی اطلاع رسانی
	<b>عنوان زیر پروژه:</b> بررسی ابعاد و تفاوت های پیکره های برجسب داده ای و پیکره های خام در زبان فارسی		
	تاریخ: 1388/04/04	ویرایش: 1/0	

### 3. پیکره های متنی برجسب داده ای و نقش آن در کاربردهای پردازش زبان طبیعی

پروژه های زیرساختی در حوزه خط و زبان فارسی را می توان پروژه هایی برای ایجاد دادگان زیرساختی در حوزه خط و زبان فارسی در نظر گرفت. از طرفی ابزارهایی نیز جهت حل مسائل زیرساختی مطرح شده در حوزه خط و زبان نیز نیاز هستند تا بتوان اقدام به تحلیل پیکره های متنی زبان جهت تولید دادگان زیرساختی نمود. بنابراین پروژه های زیرساختی خط و زبان فارسی تلفیقی از ابزارها و دادگان زیرساختی زبان خواهد بود. در این مستند، لزوم برجسب گذاری پیکره های متنی مورد بحث قرار می گیرد. بررسی ابزارهای برجسب گذاری و مقایسه ی میان آن ها، از حوزه ی مطالب این مستند خارج است. برجسب گذاری پیکره های زبانی به طور کلی می تواند در چهار سطح زبانی انجام شود که عبارتند از (بی جن خان، 1381):

- 1- تعیین برجسب مقوله کلمه<sup>1</sup>
- 2- برجسب گذاری نحوی: که شامل پردازش جملات و به دست آوردن درخت نحوی آن ها است.
- 3- برجسب گذاری معنایی: که عبارت است از استخراج صورت منطقی جملات و به دست آوردن تعبیر معنایی آن ها. پیش نیاز این نوع برجسب دهی، انجام برجسب دهی مقوله کلمات است، به این تعبیر که پیش از آن که تعبیر معنایی کلمات و جملات یک متن به دست داده شود، تعیین مقوله کلمات آن ها ضروری است.
- 4- برجسب گذاری کاربردشناختی<sup>1</sup>: که عبارت است از تعیین روابطی که میان دو کلمه در یک متن وجود دارد. به عنوان مثال مشخص کردن ضمائر و مرجع آن ها در متن در حوزه برجسب دهی کاربردشناختی قرار می گیرد.

<sup>1</sup> Wordclass tagging

	<b>عنوان پروژه:</b> فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 شورای عالی اطلاع رسانی
	<b>عنوان زیر پروژه:</b> بررسی ابعاد و تفاوت های پیکره های برجسب داده ای و پیکره های خام در زبان فارسی		
	تاریخ: 1388/04/04	ویرایش: 1/0	

با توجه به سطوح زبانی در برجسب گذاری پیکره های زبانی، ابتدا باید یک مجموعه برجسب برای برجسب گذاری یک پیکره طراحی شود. تعیین برجسب مقولات اصلی، اولین مرحله در تعیین برجسب های موجود در مجموعه برجسب است. بر این اساس در تهیه پیکره متنی زبان فارسی نیز باید مجموعه برجسب به گونه ای باشد که تمام مقوله های دستوری و انواع کلمات زبان فارسی را دربرگیرد و علاوه بر این مطابق با استانداردهای برجسب گذاری متون در زبان های دیگر مانند انگلیسی و زبان های اروپایی باشد.

برجسب هایی که برای یک پیکره در نظر گرفته می شود به طور کلی به سه دسته برجسب تقسیم می شوند که عبارتند از (Cloern, 1999):

1- برجسب های نحوی-ساختواژی<sup>2</sup>: که اصلی ترین برجسب ها هستند. این برجسب ها شامل مقوله های نحوی اصلی از جمله: فعل، اسم، صفت، قید و غیره هستند. اغلب کلماتی که در متون وجود دارند به یکی از این مقوله های اصلی تعلق دارند. اصلی ترین مقوله های نحوی که در اغلب پیکره های زبانی در نظر گرفته می شوند شامل مقوله اسم، صفت، حرف اضافه، حرف ربط، حرف تعریف، قید و عدد هستند.

2- برجسب های خاص<sup>3</sup>: که شامل کلماتی هستند که در طبقه مقوله های اصلی قرار نمی گیرند، اما تعیین برجسب آنها در استخراج اطلاعات زبانی از پیکره حائز اهمیت است. ادات شرط، حرف ندا، تکواژ صفت ساز از این دسته اند.



3- برجسب های متفرقه<sup>4</sup>: نیز شامل کلماتی است که در طبقه مقوله های اصلی قرار نمی گیرند اما در متن وجود دارند و به عنوان کلمات مجزا استخراج شده اند. کلمات خارجی، نشانه ها و علائم ریاضی و نیز علائم اختصاری در این طبقه قرار می گیرند.

<sup>1</sup> Pragmatic tagging

<sup>2</sup> Morphosyntactic

<sup>3</sup> Unique tags

<sup>4</sup> Residual /Miscellaneous tags

	<b>عنوان پروژه:</b> فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 شورای عالی اطلاع رسانی
	<b>عنوان زیر پروژه:</b> بررسی ابعاد و تفاوت های پیکره های برچسب داده ای و پیکره های خام در زبان فارسی		
	تاریخ: 1388/04/04	ویرایش: 1/0	

جدای از این تقسیم بندی برچسب های مجموعه برچسب پیکره متنی زبان فارسی یا برچسب اصلی اند یا برچسب هایی هستند که به عنوان زیربخش های برچسب های اصلی استفاده می شوند. به عنوان مثال کلمه ای مانند "کتابهایم" در پیکره به صورت زیر نمایش داده شده است:

N            N, COM,  
                  PL,1      کتابهایم

ستون اول از چپ (N) برچسب اصلی کلمه، ستون دوم (N, COM, PL,1) برچسب سلسله مراتبی<sup>1</sup> و ستون آخر (کتابهایم) خود کلمه می باشد.



برچسب گذاری ادات سخن<sup>2</sup> عمل انتساب برچسب های واژگانی به کلمات و نشانه های تشکیل دهنده یک متن است، به صورتی که این برچسب ها نشان دهنده نقش کلمات و نشانه ها در جمله باشند. درصد بالایی از کلمات از نقطه نظر برچسب واژگانی دارای ابهام هستند، زیرا کلمات در جایگاه های مختلف برچسب های واژگانی متفاوت دارند. بنابراین برچسب گذاری واژگانی عمل ابهام زدایی از برچسب ها با توجه به زمینه<sup>1</sup> مورد نظر است. برچسب گذاری واژگانی عملی اساسی برای بسیاری از حوزه های دیگر پردازش زبان طبیعی از قبیل ترجمه ماشینی، خطایاب و تبدیل متن به گفتار می باشد.

برچسب گذاری پیکره ها به عنوان یکی از مهم ترین حاشیه نویسی ها برای پیکره ها مطرح است. برچسب گذاری کاربردهای فراوانی در استخراج الگوهای اجزای واژگانی کلام در زبان دارند، به عبارت دیگر دستورات و قوانین زبان را می توان از طریق برچسب گذاری به مدلی آماری تبدیل نمود که قابل استفاده در زبان شناسی رایانه ای هستند. با استفاده از چنین مدل هایی می توان از کلمات و عبارات رفع ابهام کرد، عبارات نادرست دستوری را تشخیص داد، همچنین می توان از آن به عنوان مقدمه ای جهت استفاده از لایه های معنایی زبان یاد کرد.

در ادامه ی این بخش، نقش و ارتباط پیکره های برچسب داده ای در برخی از کاربردهای پردازش زبان طبیعی را مورد بررسی قرار خواهیم داد.

<sup>1</sup> Hierarchical tag

<sup>2</sup> Part of Speech - POS

	<b>عنوان پروژه:</b> فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 شورای عالی اطلاع رسانی
	<b>عنوان زیر پروژه:</b> بررسی ابعاد و تفاوت های پیکره های برجسب داده ای و پیکره های خام در زبان فارسی		
	تاریخ: 1388/04/04	ویرایش: 1/0	



### 3-1. خلاصه سازی

خلاصه سازی عبارت است از نمایش فشرده و دقیق متن ورودی به نحوی که متن خروجی مفاهیم مهم متن ورودی را در برداشته باشد. همگان بر این باورند که فرآیند خلاصه سازی خودکار متون باید بر پایه فهم کامل متن باشد و فرآیندی که انسان در این مورد طی می نماید را به نحوی تقلید نماید. به طور کلی خلاصه سازی یکی از زیر مجموعه مشکلات پیچیده پردازش زبان طبیعی است که هنوز بصورت کامل حل نشده است.

هدف خلاصه سازی خودکار آن است که مراحلی که توسط انسان انجام می شود، شناسایی شده و در خلاصه سازی خودکار انجام شود. بدین صورت که تمام متن خوانده و فهمیده شود و سپس خلاصه تولید شود. فهمیدن متن شامل تشخیص قسمت های مهم و غیر مهم متن است. مرحله تبدیل متن ورودی به متن خلاصه شامل مشخص کردن کلمات کلیدی، مفهوم اصلی، کلمه های مهم و جمله های مهم می باشد. در خلاصه سازی باید در نظر داشت که: (1) به چه صورت یک توصیف نرمال از متن ورودی می توانیم داشته باشیم. (2) چگونه مهمترین بخش متن را تشخیص دهیم (3) چگونه متن را تجزیه کرده و متن خلاصه را ایجاد نماییم. وجود یک سیستم تطابق دهنده که متن ورودی انتخابی را به پایگاه داده از قبل ایجاد شده مرتبط نماید ضروری به نظر می رسد.

خلاصه سازی از دید ماشین، شامل درک زبان طبیعی است. درک زبان طبیعی به اطلاعاتی همچون «دانش صوت شناسی»، «دانش مورفولوژی»، «دانش نحوی» و «دانش عملی» و «دانش سخن» نیازمند است.

برخلاف زبان انگلیسی که در آن هم حروف و هم لغات کاملاً متمایز از یکدیگرند، در زبان فارسی پیوستگی میان برخی علائم با لغات وجود دارد و علاوه بر آن تنوع نگارش در کلمات نیز موجود می باشد. ریشه یابی فعل که یکی از مراحل مهم پیش پردازش متن برای خلاصه سازی می باشد، در زبان فارسی چالش های خاص خود را دارد. به عنوان مثال در یک لغت به هم پیوسته هم بن فعل، شناسه، علامت زمان فعل و حتی شناسه های مفعولی می توان داشت که کار پردازش لغات را پیچیده تر می نماید به طوری که نمی توان از دانش، تجربه و نرم افزارهای موجود در این زمینه استفاده نمود و تولید نرم افزاری که قادر به حل تمامی این پیچیدگی ها باشد، فرایندی زمان بر و مستلزم تلاش فراوان می باشد.

	<b>عنوان پروژه:</b> فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 شورای عالی اطلاع رسانی
	<b>عنوان زیر پروژه:</b> بررسی ابعاد و تفاوت های پیکره های برجسب داده ای و پیکره های خام در زبان فارسی		
	تاریخ: 1388/04/04	ویرایش: 1/0	

مسائل و خصوصیات نگارشی زبان فارسی، امر پیش پردازش خودکار خلاصه سازی متون را دچار چالش می نماید. این چالش های شامل چالش های نگارشی و چالش های ریشه یابی است. با در دست داشتن پیکره ی برجسب داده ای، برخی از چالش های نگارشی و ریشه یابی را می توان مرتفع نمود.

خلاصه سازی متن از سه گام اصلی تشکیل شده است. مرحله پیش پردازش متن، مرحله پردازش متن و تولید خلاصه. مرحله پیش پردازش شامل حذف اطلاعات اضافی و غیر مهم و ریشه یابی می باشد. به عبارت دیگر، پیش پردازش متن در این بخش انجام می شود. بعد از پیش پردازش، بخش مهم متن باقی می ماند. بخش اصلی خلاصه سازی همان تفسیر متن و در حقیقت مرحله پردازش متن می باشد. این مرحله شامل پیدا کردن و امتیازدهی به کلمه های مهم می باشد. در مرحله آخر، خلاصه متن بر اساس جمله ها و امتیازدهی مربوطه ایجاد خواهد شد.

اهمیت تحلیل متن برای هر کاربردی در پردازش زبان طبیعی کاملاً روشن است اما اهمیت درک متن برای هر کاربردی متفاوت است. مثلاً درک سوال در یک سیستم پرسش و پاسخ برای سیستم بسیار مهم تر از تحلیل و درک متن در یک سیستم تشخیص صحبت می باشد. به هر حال مهمترین قسمت در هر سیستم کاربردی پردازش زبان طبیعی تحلیل متن ورودی خواهد بود. سطوح مختلف تحلیل متن که بیانگر هدف تحلیل می باشند را می توان به صورت ذیل تقسیم کرد:

ü تصحیح متن: تبدیل جمله به فرم استاندارد.



ü تجزیه جمله

ü استخراج اطلاعات ویژه: هدف در این حالت استخراج اطلاعات خاص می باشد

ü درک معنی متن: تبدیل هر جمله به یک فرم انتزاعی که بیانگر معنای جمله باشد. هدف از تحلیل در این مورد پی بردن به معنای جمله و یا جملات است.

ü یافتن موضوع متن: تفسیر هر قسمت یا بخش از متن برای یافتن مفهوم، موضوع یا عنوانی که متن در توضیح آن آمده است.

ü ارتباط متون: یافتن ارتباط های منطقی معنایی و ظاهری بین جملات مختلف بر اساس تشابه یا تفاوت بین جملات یا معنای آن ها.

	<b>عنوان پروژه:</b> فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 شورای عالی اطلاع رسانی
	<b>عنوان زیر پروژه:</b> بررسی ابعاد و تفاوت های پیکره های برجسب داده ای و پیکره های خام در زبان فارسی		
	تاریخ: 1388/04/04	ویرایش: 1/0	

لازم به ذکر است حسب پردازشی که روی متن انجام خواهد پذیرفت ممکن است یک و یا چندین مورد از موارد بالا در مرحله پیش پردازش متن انجام نشود.

بخشی از تحلیل متن استخراج اطلاعات است. این استخراج اطلاعات می تواند با استفاده از برجسب زنی کلمات متن باشد. برجسب زنی کمک قابل توجهی به تجزیه صحیح جملات و تصمیم گیری در مورد کلماتی که در شرایط مختلف معانی متفاوتی دارند، خواهد نمود. با تعیین دقیق نوع برجسب های مورد استفاده در پیکره های برجسب داده ای، می توان مرحله ی استخراج اطلاعات را تسهیل کرد.

تعیین کلمات کلیدی و استخراج واژه های موجود در متن از جمله دیگر کارهایی است که در بخش استخراج اطلاعات در تحلیل متن به کار می آید

## 3-2. ترجمه ماشینی

ترجمه ماشینی عبارت است از تبدیل یک متن از زبان مبدأ به زبان مقصد با استفاده از یک نرم افزار ترجمه کننده؛ چه این عمل با کمک انسان باشد چه بدون کمک انسان. نرم افزارهای مترجم، در بهترین حالت، عمل ترجمه را با دقتی در حدود 70 درصد انجام می دهند. برای به دست آوردن نتیجه بهتر، لازم است قبل و بعد از ترجمه، مقداری ویرایش روی متن انجام شود.

ترجمه ماشینی شاهد 3 نسل از ابزارها و روش ها بوده است. صرف نظر از دوره ی اول، امروزه روش های مورد استفاده در دوره ی دوم و سوم کاربرد بیشتری دارند.



رویکردهای مبتنی بر آمار و مبتنی بر مثال، از سال 1989 به بعد مورد استفاده قرار گرفته اند. ویژگی این رویکردها، قید-محوری و توجه بیشتر آن ها به جنبه لغوی در سیستم های قاعده-محور<sup>1</sup> و سیستم های تک زبانه بود.

دوره بعد از سال 1990 به بعد بوده و نسل سوم ماشینهای ترجمه در این زمان بوجود آمدند که در این نسل ماشینهای ترجمه از ترکیب رویکرد قاعده-محور نسل دوم، با روشهای متکی بر پیکره<sup>2</sup> بوجود آمده بودند. ویژگی های این نسل شامل تحلیل دستوری در سطح وابستگی ارتباطی و تحلیل معنایی در سطح

<sup>1</sup> Rule-based

<sup>2</sup> Corpus-based



	<b>عنوان پروژه:</b> فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 شورای عالی اطلاع رسانی
	<b>عنوان زیر پروژه:</b> بررسی ابعاد و تفاوت های پیکره های برجسب داده ای و پیکره های خام در زبان فارسی		
	تاریخ: 1388/04/04	ویرایش: 1/0	

تشخیص جملات و اطلاعات لغوی از فرهنگ لغات بدست آمده که با مشخصه دستوری در زبان مبدا نگهداری شده، است. متد های مبتنی بر مثال و آماری نیز به کار گرفته شدند، همچنین دانش جمع آوری شده محدود به حوزه مشخصی گردیده و از بازخورد کاربران ماشین هم در جمع آوری اطلاعات استفاده شد.

استفاده از تکنیک های مبتنی بر پیکره منجر به تلاش برای ترجمه ی متون پیشرفته تر شد. چرا که این رویکردها تفاوت های گونه شناسی زبان<sup>۱</sup>، بازشناسی عبارات، ترجمه ی اصطلاحات و هم چنین جداسازی ناپهنجاری در متون را بهتر پوشش می دهد.

با دردست داشتن داده های کافی، این تکنیک های ترجمه ی ماشینی قادر به ترجمه ی تقریبی متن از زبان مبدأ به زبان مقصد هستند. در حقیقت، مشکل اصلی در این تکنیک ها جمع آوری مقدار کافی از متن هایی با نوع خاص است. در روش های آماری پیکره های چندزبانیه ی عظیمی مورد نیاز است، در حالی که در روش های مبتنی بر نحو چنین نیست. اما در عین حال، در توسعه ی روش های مبتنی بر نحو به زبان شناس های متبحر نیاز است تا بتوان گرامر زبان را بدرستی و با دقت بالایی طراحی کرد. مشکل دیگر روش های مبتنی بر نحو، دقیق نبودن گرامرهای فعلی برای برخی از زبان ها است. زبان فارسی از جمله ی این زبان ها است.

بنابر توضیحات فوق، امروزه 2 رویکرد مختلف در فن آوری های ترجمه ماشینی را می توان شناسایی کرد: ترجمه ی مبتنی بر قاعده<sup>۲</sup>، و ترجمه ی مبتنی بر پیکره<sup>۳</sup>.



رویکردهای مبتنی بر پیکره، از داده های موجود در پیکره های متنی استفاده می کنند.

رویکرد اول این نوع از ماشین های ترجمه، مبتنی بر مثال نام دارند که در آنها در پیکره یا مجموعه اطلاعات زبانی به دنبال مشابه و نمونه ای که قبلا آمده بوده و اینکه چگونه ترجمه شده بوده می گردد.

<sup>۱</sup> Linguistic typology

<sup>۲</sup> Rule-Based MT

<sup>۳</sup> Corpus-Based MT

	<b>عنوان پروژه:</b> فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 شورای عالی اطلاع رسانی
	<b>عنوان زیر پروژه:</b> بررسی ابعاد و تفاوت های پیکره های برجسب داده ای و پیکره های خام در زبان فارسی		
	تاریخ: 1388/04/04	ویرایش: 1/0	

رویکرد آماری، زبان مبدا را بدون استفاده از قوانین به زبان مقصد نگاشت می کند و از مدل زبانی و اطلاعاتی که در پیکره وجود دارد استفاده می کند. به طور کل دانش زبانی کمی در این روش بکار گرفته می شود و از احتمال رخ دادن ترکیبات جملات استفاده می نماید.

مراحل یک ترجمه به صورت کلی عبارت است از :



1. تفسیر پیکره زبان مقصد به صورت غیربرخط
2. تفسیر جملات زبان مبدأ
3. بهره برداری از پیکره زبان مقصد برای ارائه بهترین ترجمه
4. ترکیب و ساخت ترجمه خروجی

در مرحله ی تفسیر پیکره زبان مقصد، عبارت ها اندیس گذاری شده و طبقه بندی می شود. این عمل باعث می شود که جستجوها برای یافتن بهترین تطبیق، با سرعت بیشتری انجام شود.

تفسیر جملات زبان مقصد درست قبل از اینکه به الگوریتم تطبیق داده شوند انجام میگیرد. بدین ترتیب اطلاعات زبانی لازم در اختیار این الگوریتم قرار می گیرد. ابتدا، برجسب گذاری انجام شده و توسط ابزارهای بخش کننده<sup>1</sup> به قسمت های متنهایی بخش بندی شده می شود. سپس به قسمت های برجسب دار، مجموعه لغات دوزبانه اعمال شده و ابتدا ترجمه روی ترکیبات کلمه ای و سپس روی تک تک کلمات صورت می گیرد. خروجی این مرحله، کلیه ی ترجمه های ممکن به زبان مقصد هستند. بعد از این مرحله، از این منابع برای ساخت مجموعه دو زبانی لغات استفاده شده و پس از بررسی سنخیت و دقت آن با نیاز سیستم، همگون می گردد.

در برخی از روش ها نیز پیشنهاد شده است که می توان با استفاده از یک پیکره ی متنی موازی و توصیف گرامری مجموعه ی کوچکی از جملات پیکره (بصورت دستی)، گرامر یک زبان را یاد گرفت.

از چالش های ترجمه ی ماشینی در روشهای ترکیبی (مبتنی بر قوانین و آماری)، می توان به «فراهم کردن دادگان آن» و ابزار تجزیه کننده اشاره کرد. کلیه ی موارد زیرشاخه ی فراهم کردن دادگان، در روش های و نیز برخی از موارد زیرشاخه ی ابزار تجزیه کننده، در مراحل تهیه ی پیکره های برجسب داده ای (یک زبانی و چندزبانی) مورد پوشش قرار می گیرد.

	<b>عنوان پروژه:</b> فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 شورای عالی اطلاع رسانی
	<b>عنوان زیر پروژه:</b> بررسی ابعاد و تفاوت های پیکره های برچسب داده ای و پیکره های خام در زبان فارسی		
	تاریخ: 1388/04/04	ویرایش: 1/0	

## 4. نتیجه گیری

در بخش پیش، ارتباط و نیازمندی برخی کاربردهای پردازش زبان طبیعی به پیکره‌ی متنی را برشمردیم. همان‌گونه که آمد، پیکره‌های برچسب داده‌ای، حاوی اطلاعات ارزشمندی هستند که تقریباً در تمامی کاربردهای پردازش زبان طبیعی می‌تواند مورد استفاده قرار گیرد. در این مستند تنها به دو کاربرد از کاربردهای پردازش زبان طبیعی پرداختیم، لکن امتیازات پیکره‌های برچسب داده‌ای، در دیگر کاربردها نیز مشاهده می‌شود. لازم بذکر است که برتری پیکره‌های برچسب داده‌ای به پیکره‌های خام، منحصر به زبان خاصی نیست و به تمامی زبان‌ها قابل تعمیم است.