

## گزارشی از ساخت نخستین پیکره چند زبانه برای زبان فارسی

بهرنگ قاسمی‌زاده\*، سعید رحیمی، مرتضی سالاریان، علی بهاری سلیم

### چکیده

این مقاله، اولین پیکره فارسی موازی با تعداد زیادی از زبان‌های اروپایی را معرفی می‌نماید. این مقاله، اولین قدم‌ها را برای ساخت منابع اساسی جهت پردازش زبان فارسی معرفی می‌نماید. این مرحله از کار شامل معرفی ویژگی‌های صرفی-نحوی زبان فارسی و رمزنگاری آنها بر پایه مدل EAGLES/MULTEXT و منابع خاص MULTEXT-East می‌باشد. این مقاله پس از معرفی مختصر زبان فارسی، با تأکید بر شیوه نگارش و ویژگی‌های صرفی-نحوی آن، به ارائه یک دسته‌بندی جدید برای مقوله‌های گفتاری فارسی پرداخته و رسم الخطی را جهت ارائه فارسی در محیط دیجیتال ارائه می‌نماید. پس از آن آماری از پیکره ساخته شده ارائه می‌شود. کار صورت گرفته منحصرًا توسط افراد داوطلب و بدون هیچگونه پشتیبانی مالی و یا معنوی از سازمان خاص صورت گرفته است.

### واژه‌های کلیدی

پیکره، زبانشناسی رایانه‌ای، ویژگی‌های صرفی-نحوی، فارسی، حاشیه نویسی پیکره‌ها، پردازش زبان طبیعی.

## Report of Constructing the First Parallel, Multilingual Corpus of Farsi

Behrang Qasemizadeh, Saeed Rahimi, Morteza Salarian, Ali Bahari Salim  
Text and Speech Technology Corporation

### Abstract

In this paper, we have introduced the first parallel corpus of Farsi with more than 10 other European languages. This paper describes primary steps toward preparing Basic Language Resources Kit (BLARK) for Farsi. Since now, we have proposed morphosyntactic features of Farsi based on EAGLE/MULTEXT models and specific resources of MULTEXT-East. The paper introduces Farsi Language, with emphasize on its writing system and morphosyntactic features, then a new Part-of-Speech categorization and orthography for Farsi in digital environments is proposed. Finally, the corpus and related static will be described. This work has been done in an informal way and with a volunteer based team.

### Keywords

Natural Language Processing, Morphosyntactic Features, Orthography, Corpus, Computational Linguistics, Corpus Annotation, Parallel Corpora, Farsi

\* سرپرست تیم پدیدآورندگان پیکره، کارشناس ارشد هوش مصنوعی، شرکت فن‌آوری متن و گفتار، [qasemizadeh@comp.iust.ac.ir](mailto:qasemizadeh@comp.iust.ac.ir)

## ۱- مقدمه

امروزه با افزایش اهمیت فن آوری ارتباطات و اطلاعات، نیاز به فن‌آوری‌های زبان و گفتار نیز افزایش یافته است. برای اینکه افراد بتوانند از زبان مادری خود بر روی رایانه‌ها استفاده نمایند، نیاز به یک مجموعه تدارکات اولیه (همانند واژگان، پیکره و مجموعه‌ای از ابزارها) است. تلاش‌های زیادی جهت مهیا نمودن بسته منابع زبانی پایه<sup>۱</sup> برای پردازش زبانها، به خصوص زبان‌هایی که از لحاظ تجاری کمتر مورد توجه اند صورت گرفته است. حاصل این تلاش‌ها، فراهم آمدن پیکره‌ها، ابزارهای تحلیل زبان همانند تحلیلگر صرفی و ... است. علاوه بر این، این منابع به راحتی قابل دسترسی و برای مقاصد آموزش و تحقیق کاملاً رایگان است. متأسفانه برای زبان فارسی تا به حال کارهای مختصری در این زمینه صورت گرفته است. [۱][۲] علاوه بر آنکه کارهای صورت گرفته نیز مطابق با استاندارد و یا چارچوب کاری خاص صورت نگرفته است و اکثر آنها جهت استفاده محققین و دانشجویان در دسترس نمی‌باشد.

این مقاله به اولین قدم‌های صورت گرفته جهت فراهم آوردن منابع اولیه جهت پردازش زبان فارسی اشاره دارد. این شامل معرفی رسم الخط (قراردادی جهت ارائه پیکره) برای فارسی در محیط دیجیتال و رمزنگاری آن، معرفی مقوله‌های گفتاری، و معرفی ویژگی‌ها صرفی-نحوی فارسی منطبق بر مدل‌های ارائه شده توسط EAGLE/MULTEXT و استاندارد ارائه شده در MULTEXT-East است. [۳]

پروژه MULTEXT-East گسترشی از پروژه EU-MULTEXT بود که به گسترش منابع زبانی برای شش زبان شامل بلغاری، چک، استونی، لهستانی، رومانیایی، اسلونی و زبان انگلیسی به عنوان زبان مرکزی<sup>۲</sup> در پروژه پرداخت. نتیجه اصلی حاصل از این پروژه یک پیکره چند زبانه حاشیه نویسی شده، و مجموعه‌ای از ابزارها برای این هفت زبان است. این پروژه برای دیگر زبان‌های اروپایی نیز در حال گسترش و پیاده سازی است. کارآمدترین قسمت از این پروژه شامل منابع تشریح ویژگی‌های صرفی-نحوی (MSD)<sup>۳</sup> است که از سه لایه تشکیل شده است. این سه لایه به ترتیب انتزاع به شرح ذیل است [۵]:

- MSD 1984: پیکره ۱۹۸۴ (کتاب ۱۹۸۴ اثر جرج ارول) که با ویژگی‌های صرفی-نحوی حاشیه نویسی شده است. در این پیکره هر یک از واژه‌ها با ویژگی‌های صرفی-نحوی رفع ابهام شده با توجه محل ظهور آنها، حاشیه نویسی شده اند.
- MSD Lexicons: واژگانی از واژه‌های صرف شده که در پیکره ظاهر شده اند به همراه لما و ویژگی‌های صرفی-نحوی آنها.
- MSD Specs: شرح ویژگی‌های صرفی-نحوی. این شرح مشخص می‌نماید که چه MSD‌هایی برای یک زبان معتبر است و معنای آن چیست برای مثال Nems به معنی PoS اسم، نوع اسم عام، جنس خنثی، و عدد برار مفرد است.

بدین ترتیب MULTEXT-East چارچوبی قابل فهم برای گسترش پیکره‌ها فراهم می‌آورد. علاوه بر این منابع مختلفی مبتنی بر این چارچوب وجود دارد. به عنوان مثال پیکره ۱۹۸۴ برای چندین زبان مختلف وجود دارد. نکته جالب توجه اینکه ترجمه رمان ۱۹۸۴ به فارسی وجود دارد و از سوی دیگر با توجه به ماهیت زبان فارسی به عنوان یک زبان هندی-اروپایی امکان اضافه کردن آن به این چارچوب امکان پذیر است. بنا بر آنچه که گفته شد و با توجه به نیازهای موجود برای پردازش رایانه‌ای فارسی اولین قدم‌ها جهت اضافه نمودن فارسی به این چارچوب برداشته شد. در کار گذشته نویسندگان این مقاله [۳] فاکتورهای مهم در این زمینه، شامل رسم الخط فارسی و نحوه رمزنگاری آن در محیط الکترونیکی، زبان فارسی و ویژگی‌های صرفی و نحوی آن، و در نهایت تعامل میان این فاکتورها و چگونگی تاثیر آنها بر فرایند پردازش زبان به تفصیل بحث شده است.

در این مقاله، ادامه فعالیت‌های صورت گرفته در این جهت، یعنی ساخت پیکره فارسی و استفاده از تئوری‌های ارائه شده در [۳] در یک کاربرد واقعی گزارش شده است. با توجه به اهمیت موضوع در بخش ۲ ابتدا خلاصه‌ای از کارهای صورت گرفته ارائه می‌شود. نسخه جدید از آخرین تغییرات اعمال شده بر ویژگی‌های صرفی-نحوی و جداول مرتبط با آن مطالب بخش ۳ را تشکیل می‌دهند. بخش ۴ به ارائه پیکره ساخته شده و آمار مرتبط با آن می‌پردازد. بخش ۵ از مقاله به نتیجه گیری و شرح کارهای آتی اختصاص داده شده است.

## ۲- زبان فارسی در چارچوب MULTEXT-East

همانطور که گفته شد، در کار قبلی، زبان فارسی در چارچوب MULTEXT-East گنجانده و مطابق با آن ویژگی‌های صرفی-نحوی زبان فارسی ارائه شده است. در پاره‌ای موارد نیز با توجه به برخی ویژگی‌های خاص زبان فارسی، گروهی از ویژگی‌ها به چارچوب ذکر شده اضافه شده است. اضافه کردن فارسی به این چارچوب با چالش‌های عمده‌ای روبرو بود مهمترین چالش‌ها از این میان شامل موارد ذیل است:

- مشکلات مرتبط با رمزنگاری و ارائه فارسی در محیط الکترونیکی، و همزمان با آن رسم الخط فارسی در محیط الکترونیکی
  - تعریف ویژگی‌های صرفی-نحوی مطابق با واقعیات زبان فارسی و در نظر گرفتن کاربرد تعاریف ارائه شده در زبانشناسی رایانه‌ای و تحلیل رایانه‌ای زبان فارسی
  - تعامل و اثر متقابل تعاریف ارائه شده در بند ۱ و ۲، و همزمان با آن تاثیر بر فرآیندهای تحلیل همانند توکن بندی متون الکترونیکی فارسی و دیگر تحلیل‌های بعدی. (شکل ۱)
- برای حل مشکل اول یعنی رسم الخط و رمزنگاری متون الکترونیکی، پس از مطالعه و مشورت با اساتید امر، تصمیم بر آن شد تا روش ارائه متون الکترونیکی فارسی ترکیبی از استاندارد ارائه شده توسط موسسه ملی استاندارد برای رمزنگاری کاراکترهای فارسی در

Unicode [۶] و همزمان با آن رسم الخط رسمی فارسی [۷] ارائه شده توسط فرهنگستان ادب و زبان فارسی برای سیستم مبتنی بر کاغذ (paper-based) باشد. در حقیقت روش ارائه شده برای رسم الخط الکترونیکی فارسی تطابق [۷] با [۶] است. مطابق روش ارائه شده، کاراکتر کنترلی ZWNJ به جای نیم‌فاصله در متون ظاهر می‌شود و مطابق آن کلیه پیش‌وندها و پسوندها، در زمانیکه مطابق رسم الخط باید جدا از ریشه کلمه ظاهر شوند، به همراه یک ZWNJ در ابتدای پسوند و یا در انتهای پیشوند رمزنگاری و نمایش داده می‌شوند. کاراکتر فاصله تنها در میان دو کلمه مستقل فارسی استفاده می‌شود. شرح کامل تر از این مطلب در [۳] و [۸] آمده است.

جدول (۱). مقوله های گفتاری در فارسی به همراه تعداد ویژگی ها

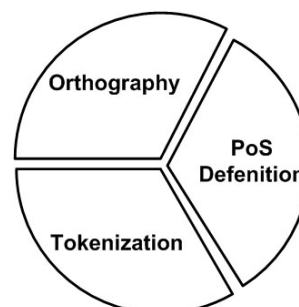
Part of Speech	Code	Number of Attributes
Noun	N	4
Verb	V	10
Adjective	A	4
Pronoun	P	6
Determiner	D	1
Adverb	R	2
Adposition	S	2
Conjunction	C	2
Numeral	M	3
Interjection	I	0
Abbreviation	Y	0

### ۳- ویژگیهای صرفی-نحوی فارسی

روش ارائه ویژگی ها در چارچوب MULTEXT-East، روشی مبتنی بر مکان است و هر ویژگی تنها با یک کاراکتر نمایش داده می‌شود. در این چارچوب اولین کاراکتر از یک حاشیه نویسی مقوله گفتاری واژه را مشخص می‌نماید. سپس مقادیر ویژگی ها با توجه به محل از پیش تعریف شده برای آن، در جلوی این کاراکتر نوشته می‌شوند. در صورتیکه یک ویژگی برای یک زبان قابل تعریف نباشد، در مکان مرتبط با این ویژگی خط تیره (-) قرار می‌گیرد. جهت اطلاع بیشتر از روش ارائه ویژگی ها در این چارچوب می‌توان به [۱۴] مراجعه نمود.

در جداول پیش رو ویژگی های صرفی-نحوی فارسی برای هر مقوله آورده شده است. اولین ستون از سمت چپ در این جداول محل قرارگیری مقادیر آن ویژگی را در یک حاشیه نویسی نمایش می‌دهد. برای کوتاه نویسی، ویژگی های غیر قابل کاربرد برای فارسی از جداول حذف شده اند. بنابراین ممکن است در برخی جداول مقادیر برای برخی مکان ها تعریف نشده باشد، در این صورت در محل این ویژگی‌ها همواره یک خط تیره قرار می‌گیرد. ستون دوم در این جدول نام ویژگی را نمایش می‌دهد به عنوان مثال ویژگی شخص برای مقوله گفتاری فعل. ستون سوم و چهارم به ترتیب مجموعه مقادیر با معنی برای این ویژگی و کاراکتر نمایش دهنده این مقدار را نمایش می‌دهند. به این ترتیب برای ویژگی شخص برای مقوله گفتاری فعل، سه مقدار با معنی شامل اول، دوم و سوم وجود دارد که به ترتیب با کاراکترهای ۱، ۲ و ۳ در مکان چهارم پس از کاراکتر نمایش دهنده مقوله گفتاری فعل (V) قرار می‌گیرند.

جدول ۲ تا ۱۰ ویژگی های صرفی-نحوی برای مقولات گفتاری مختلف، نحوه رمزنگاری، و مقادیر در نظر گرفته شده برای آنها را نمایش می‌دهد. این جداول برای دو مقوله گفتاری صوت و اختصارات، به دلیل در نظر گرفتن هیچ ویژگی صرفی-نحوی نیامده است.



شکل ۱. اگر تمامی شکل را به عنوان استاندارد ارائه شده برای برچسب دهی پیکره در نظر بگیرید، آنگاه سیاست اتخاذ شده در توکن بندی، دسته بندی ارائه شده برای مقوله های گفتاری، و رسم الخط زبان، المانهای اساسی هستند که مستقیماً بر روی مجموعه برچسب های تعریف شده برای حاشیه نویسی پیکره تاثیر می‌گذارند.

با مطالعات صورت گرفته بر روی زبان فارسی [۹][۱۰][۱۱][۱۲][۱۳] تقسیم بندی جدیدی از مقوله های فارسی جهت ارائه آن در چارچوب MULTEXT-East ارائه شده است. این دسته بندی جدید قائل به یازده مقوله گفتاری برای زبان فارسی است. جدول شماره یک این مقوله ها را به همراه تعداد ویژگی های در نظر گرفته شده برای آنها نمایش می‌دهد؛ برای هر یک از این مقوله‌ها ویژگی های صرفی-نحوی خاص در نظر گرفته شده است و در چارچوب MULTEXT-East ارائه شده است. به این دلیل پاره ای از ویژگی ها و یا مقادیر آنها صرفاً جهت فارسی به این چارچوب افزوده شده است. همانطور که برخی از ویژگی های از پیش تعریف شده برای دیگر زبان ها برای فارسی قابل کاربرد نیستند.

نظر نویسندگان مقاله و نتایج حاصل، بر این است که روش ارائه شده به همراه تعاریف جدید ارائه شده از مقوله های گفتاری فارسی و ویژگی های صرفی-نحوی آن بهترین نتیجه را در تحلیل رایانه ای فارسی به همراه دارد. در ادامه و در بخش بعد ویژگی های صرفی-نحوی ارائه شده برای فارسی و روش رمزنگاری آنها در چارچوب MULTEXT-East توضیح داده شده است.

جدول (۴). ویژگیهای صرفی-نحوی مقوله صفت

P	ATT	VAL	C
1	Type	qualificative	f
2	Degree	positive comparative superlative	p c s
5	Case	genitive	g
6	Definiteness	no yes	n y

جدول (۲). ویژگیهای صرفی-نحوی مقوله اسم

P	ATT	VAL	C
1	Type	common proper	c p
3	Number	singular plural	s p
4	Case	genitive	g
5	Definiteness	no yes	n y

جدول (۵). ویژگیهای صرفی-نحوی مقوله ضمیر

P	ATT	VAL	C
1	Type	personal demonstrative indefinite interrogative reflexive reciprocal	p d i q x y
2	Person	first second third	1 2 3
4	Number	singular plural	s p
5	Case	genitive accusative	g a
8	Clitic	no yes	n y

جدول (۳). ویژگیهای صرفی-نحوی مقوله فعل

P	ATT	VAL	C
1	Type	main auxiliary modal copula light	m a o c l
2	VForm	indicative subjunctive imperative participle	i s m p
3	Tense	present past	p s
4	Person	first second third	1 2 3
5	Number	singular plural	s p
8	Negative	no yes	n y
10	Clitic	no yes	n y
14	Aspect	progressive	p
15	Courtesy	no yes	n y
16	Transitive	no yes	n y

جدول (۶). ویژگیهای صرفی-نحوی مقوله Determiner

P	ATT	VAL	C
1	Type	demonstrative indefinite interrogative exclamative article exceptional	d i q e a x
4	Number	singular plural	s p

جدول (۱۱). تعداد برچسب‌های با معنی برای هر مقوله گفتاری

Part of Speech	Number of Attributes
Noun	12
Verb	639
Adjective	12
Pronoun	78
Determiner	6
Adverb	3
Adposition	3
Conjunction	4
Numeral	12
Interjection	1
Abbreviation	1
<b>Total Number</b>	<b>771</b>

#### ۴-۴. پیکره

همانطور که گفته شد، در چارچوب MULTEXT-East رمان ۱۹۸۴ اثر جرج ارول به عنوان متن اصلی و پیکره متنی انتخاب شده است. بنابراین نسخه فارسی این کتاب نیز برای اضافه کردن زبان فارسی به این چارچوب حاشیه نویسی شده است. بطور کلی این پیکره از:

- ۱۱۰۰۰۰ توکن
- ۱۱۲۶۶ پاراگراف
- ۶۶۰۶ جمله
- ۶۶۳۲ لما
- ۱۳۵۹۷ نوع کلمه

تشکیل شده است. از میان ۷۷۱ برچسب مختلف با معنی و ممکن برای فارسی تنها ۴۴۸ برچسب مختلف در این پیکره رخ داده است. این پیکره در قالب بسته ارائه شده در MULTEXT-East ارائه، و به صورت بر خط جهت تحقیق، به صورت رایگان در اختیار محققین و دانشجویان علاقمند به پردازش زبان طبیعی، قرار خواهد گرفت.

#### ۵- نتیجه

این مقاله به معرفی اولین قدم‌ها جهت ساخت منابع اساسی جهت پردازش زبان فارسی می‌پردازد. در اولین قدم جهت حرکت به سوی این هدف، ویژگی‌های صرفی و نحوی برای زبان فارسی مطابق با خط مشی‌های ارائه شده توسط EAGLE/MULTEXT ارائه شده است. کار انجام شده، دسته‌بندی جدیدی از مقوله‌های گفتاری زبان فارسی به همراه ویژگی‌های صرفی-نحوی جدید، معرفی می‌نماید. علاوه بر این، مقاله ارائه شده، روش جدید جهت ارائه متون الکترونیکی فارسی ارائه می‌نماید. به طور کلی روش انتخاب شده جهت ارائه متون و همچنین دسته‌بندی ارائه شده از مقوله‌های گفتاری در جهتی صورت گرفته است که نخست همگی ویژگی‌های صرفی-نحوی زبان فارسی را

جدول (۷). ویژگی‌های صرفی-نحوی مقوله قید

P ATT	VAL	C
2 Degree	positive comparative	p c
7 Case	genitive	g

جدول (۸). ویژگی‌های صرفی-نحوی مقوله Adpositions

P ATT	VAL	C
1 Type	preposition postposition	p t
2 Formation	simple compound	s c

جدول (۹). ویژگی‌های صرفی-نحوی مقوله حرف ربط

P ATT	VAL	C
1 Type	coordinating subordinating	c s
2 Formation	simple compound	s c

جدول (۱۰). ویژگی‌های صرفی-نحوی مقوله عدد

P ATT	VAL	C
1 Type	cardinal ordinal fractal ordinal2	c o f r
4 Case	genitive	g
6 Definiteness	no yes	n y

شرح دقیق از ویژگی‌های ارائه شده و دلیل ارائه هر یک از این ویژگی‌ها در [۳] آمده است. جدول ۱۱ تعداد برچسب‌های مختلف با معنی برای هر مقوله گفتاری را نمایش می‌دهد. جمعاً امکان وقوع ۷۷۱ برچسب مختلف با معنی برای کلمات فارسی وجود دارد. از بین مقوله‌های گفتاری مختلف فعل‌ها بیشترین تعداد برچسب را به خود اختصاص داده‌اند.

- [10] Meshkatodini M.: *Introduction to Persian Transformational Syntax*, 2nd Edition, ISBN: 964-6335-80-2, Ferdowsi University Press, (2003).
- [11] Lazard, G.: *A Grammar of Contemporary Persian*, Mazda Publishers, (1992).
- [12] Riazati D.: *Computational Analysis of Persian Morphology*, Msc Thesis, Department Of Computer Science, RMIT, (1997).
- [13] Bateni, M.: *Towsif-E Sakhteman-E Dastury-E Zaban-E Farsi [Description Of The Linguistic Structure Of Persian Language]*, Amir Kabir Publishers, Tehran, Iran, (1995).
- [14] Erjavec T.: *MULTEXT-East Morphosyntactic Specifications*, Version 3.0. Supported By EU Projects Multext-East, Concede And TELRI, (2004).

## زیرنویس‌ها

- <sup>1</sup> Basic language resources kits (BLARK)
- <sup>2</sup> The Hub language
- <sup>3</sup> Morphosyntactic Descriptions
- <sup>4</sup> Tomaz Erjavec
- <sup>5</sup> Damir Cavar

کاملاً و به شکل واقعی ارائه نماید و در عین حال حداکثر هم خوانی و سازگاری با دیگر زبانها را داشته باشد.

در قدم بعدی جهت آماده سازی منابع اساسی برای فارسی، پیکره ۱۹۸۴ به فارسی، بر چسب دهی شده است. وجود این پیکره برای زبان فارسی علاوه بر فراهم آوردن یک مجموعه داده استاندارد جهت استفاده محققین پردازش زبان طبیعی در فارسی، امکان استفاده از متدهای آماری جهت پردازش فارسی را معرفی می نماید. علاوه بر این چند زبانه بودن این پیکره، سبب میشود تا بتوان از آن در تحقیقات مرتبط با ترجمه ماشینی از فارسی/به فارسی از دیگر زبانهای موجود در این چارچوب، و در نهایت تحقیق و گسترش چنین سامانه هایی پرداخت. علاوه بر اینکه وجود زبان فارسی در یک چارچوب استاندارد و موازی با دیگر زبانها، امکان مطالعه و تحقیق جهت استخراج دانش به کمک پردازش های بین زبانی و یا حتی مقایسه بین این زبان ها و زبان فارسی، فراهم می آورد.

به عنوان کار آینده، تصمیم بر گسترش پیکره به ۱ میلیون کلمه می باشد. علاوه بر اینکه گسترش سیستمی جهت استاندارد سازی خودکار متون فارسی نیز مد نظر قرار گرفته است.

## سپاسگزاری

نویسندگان این مقاله مایلند از پرفسور توماژ اریاوتز<sup>۴</sup> و پرفسور دمیر چوار<sup>۵</sup> به دلیل هم‌فکری و راهنمایی های بسیارشان تشکر نمایند.

## مراجع

- [1] Strik, H. Daelemans, W. Binnenpoorte, D. Sturm, J. de Vriend, F. and Cucchiarini, C.: *Dutch HLT resources: From BLARK to priority lists*, In Proceedings of ICSLP, Denver, USA, pp. 1549-1552, Denver, USA, (2002).
- [2] Krauwer, S. Maegaard, B. Choukri, K. and Damsgaard Jørgensen, L.: *Report on BLARK for Arabic*, (2004).
- [3] QasemiZadeh B., Rahimi S.: *Persian in MULTEXT-East Framework*. FinTAL 2006: 541-551, Springer Publisher, Lecture Notes in Computer Science, Vol. 4139, 2006.
- [4] Ide N. and Veronis J.: *Multext: Multilingual Text Tools And Corpora*. In 15th Int. Conference On Computational Linguistics, Pages 588-592, Kyoto, Japan, (1994).
- [5] Erjavec T., Krstev C., Petkevic V., Simov K., Tadic M., and Vitas D.: *The MULTEXT-East Morphosyntactic Specifications For Slavic Languages*, Proceedings Of The EACL 2003 Workshop On The Morphological Processing Of Slavic Languages, (2003).
- [6] Isiri 6219:2002: *Information Technology – Persian Information Interchange and Display Mechanism, Using Unicode*, (2002).
- [7] Iran's Academy Of Persian Language and Literature: *Official Persian Orthography*, ISBN: 964-7531-13-3, 3<sup>rd</sup> Edition, (2005).
- [8] QasemiZadeh B., Rahimi S., Safae M., *Challenges in Persian Electronic Texts Analysis*, InScit2006, Merida, Spain, 2006.
- [9] Hasan A. and Ahmadi Givi H.: *Persian Grammar*, ISBN964-318-007-7, 22nd Edition, Tehran, (2002).