



FACULTY OF HUMAN AND SOCIAL SCIENCES

An Electronic Dictionary of Persian Verbs

Bahareh Kakanaeeni

Dissertation
Master of Language Sciences

Supervisor Prof. Doutor Jorge Baptista

Faro, 2014

Statement

An Electronic Dictionary of Persian Verbs

Declaração de autoria do trabalho

Declaro ser a autora deste trabalho, que é original e inédito. Autores e trabalhos consultados estão devidamente citados no texto e constam da listagem de referências incluída.

©2014, Bahareh Kakanaeeni/ Universidade do Algarve

A Universidade do Algarve tem o direito, perpétuo e sem limites geográficos, de arquivar e publicar este trabalho através de exemplares impressos reproduzidos em papel ou de forma digital, ou por qualquer outro meio conhecido ou que venha a ser inventado, de o divulgar através de repositórios científicos e de admitir a sua cópia e distribuição com objetivos educacionais ou de investigação, não comerciais, desde que seja dado crédito ao autor e editor.

An Electronic Dictionary of Persian Verbs

Declaration of authorship of the work

I hereby declare to be the author of this work, which is original and unpublished. Writers and works consulted are properly quoted in the text and listed in the list of references here included.

©2014 Bahareh Kakanaeeni / University of Algarve

The University of Algarve is entitled, perpetually and without any geographical limits, to file and publicizing this work through printed copies reproduced on paper or in digital form, or by any other known medium or any as yet to be invented, through its promotion on scientific reposition and admit its copy and distribution for educational or research non-commercial purposes, as long as credit is given to the author and to the publisher.

Acknowledgements

Foremost, I would like to express my honest appreciation to my supervisor Prof. Jorge Baptista for support of my Master study and thesis, for his patience, motivation, enthusiasm, and enormous knowledge. His supervision helped me in all the time of research and writing of this thesis. I could not have imagined having a better supervisor and teacher for my Master study.

My deepest heartfelt appreciation goes to my dear husband, Vahid, for his passion, kindness and sustenance, without him this project would simply stay as a nice idea.

Last but not the least, I express my love and gratitude to my families, for their understanding and endless love, through the duration of my studies.

Faro, July 3rd 2014

Bahareh Kakanaeeni

Abstract

There are more than 110 million Persian speakers in the world, but the lexical resources for Natural Language Processing (NLP) of the Persian language are still too scarce. Though Persian (or Fārsi) uses an adapted form of the Arabic script, it is an Indo-European language and its verbal inflection is based on stems and affixes.

This project is to be considered as a first step towards the construction of large-scale lexical resources for Persian, and the development of a Persian module to be distributed with the Unitex linguistic development platform. The specific goal of this dissertation is to build a morphologic, machine-readable dictionary of Persian verbs, using a dictionary of lemmas and a set of morphologic finite-state transducers (FST) to generate all the inflected forms associated to each lemma, and encode them with all the relevant morphosyntactic information (tense, person-number, etc.).

This task is complicated in Persian verbal morphology by the fact that each verb has two stems (past and present), different inflection paradigms are used for written (formal) and oral (informal) language uses, and several compound tenses can be formed through combining prefixes and suffixes with base inflected forms, the same tense being able to constitute one, two or more different tokens (separate written forms).

A small dictionary of lemmas (145) and their respective written (292) and spoken (127) stems was built, each stem was provided with the appropriate inflection conventional code, which correspond to an inflection paradigm. The list of lemmas was compiled based on frequency data from a large Persian corpus, the TEP (Tehran English-Persian Parallel Corpus), containing around 4.5 million words, and selecting the most frequent verb forms.

At its current state, the dictionary of inflected forms contains 1,536 entries.

In Persian, in average, each verb lemma yields 28 simple inflected forms, 14 simple 'written' inflected forms and 14 'spoken' inflected forms, and relatively, for each stem there are 7 inflected forms.

The recognition of compound tenses is carried out by a set of 22 FST using the system morphological mode and the previous lexical annotation of simple verb forms. These FSTs allowed for the retrieval of 3,953 compound verb tenses from the corpus.

For the evaluation of this language resource, a sample text, retrieved from the www.persian.euronews.com, and containing around 1000 words was used to assess the lexical coverage of the simple words' dictionary and the compound tenses' lexical graphs. The evaluation was done manually, based on the recognized words and those that were not identified by lexical resources built here.

Keywords

Persian (Farsi), verb, electronic (or machine-readable) dictionary, inflectional morphology, finite-state automata and transducers

Resumo

Há mais de 110 milhões de falantes da língua Persa (ou Farsi) no mundo, mas os recursos lexicais para Processamento de Língua Natural (NLP) disponíveis para esta língua ainda são muito escassos. Embora o Persa use uma forma adaptada do alfabeto árabe, é uma língua indo europeia e sua flexão verbal é baseada em morfemas radicais e afixos.

Este projeto deve ser considerado como um primeiro passo para a construção de recursos lexicais em larga escala para o Persa, e o desenvolvimento de um módulo do Persa a ser distribuído com a plataforma de desenvolvimento linguístico UNITEX. O objetivo específico deste trabalho é construir um dicionário eletrónico, legível por máquinas, de verbos do persa, usando um dicionário de lemas e um conjunto de transdutores de estados finitos (FST) morfológicos para gerar todas as formas flexionadas associadas a cada lema, codificando-as com toda a informação morfossintática relevante (pessoa-número, tempo-modo, etc.). Esta tarefa é complicada na morfologia verbal persa pelo facto de cada verbo tem dois estemas (ing. *stem*), um para os tempos do passado e outro para os do presente; são também empregues diferentes paradigmas de flexão na linguagem escrita (mais formal) e na oralidade (informal); e, finalmente, pelo facto de vários tempos compostos poderem ser formados através da combinação a prefixos e sufixos com formas flexionadas de base, ao mesmo tempo que podem constituir-se num único, em dois ou até mais diferentes palavras gráficas (*tokens*).

Um pequeno dicionário de lemas (145) e seus respectivos estemas da forma 'escrita' (292) e 'falada' (127) foi, então, construído, e cada estema acompanhado pelo código de flexão convencional de flexão adequado, e que corresponde a um paradigma de inflexão. A lista dos lemas foi compilada com base em dados de frequência de um grande corpus persa, o TEP (Tehran English –Persian Parallel Corpus), contendo cerca de 4,5 milhões de palavras, e selecionando as formas verbais mais frequentes. No seu estado atual, o dicionário de formas flexionadas contém 1.536 entradas. Em Persa, em média, cada lema verbal produz 28 formes simples: 14 formas simples 'escritas' flexionadas e 14 formas 'faladas', dado que cada stema.

O reconhecimento de tempos compostos é realizado por um conjunto de 22 FSTs utilizando o modo morfológico do sistema e a prévia anotação lexical de formas verbais simples. Estes grafos permitiram reconhecer 3.953 formas compostas no corpus utilizado.

Para a avaliação deste recurso linguístico, utilizou-se um texto obtido a partir da www.persian.euronews.com e contendo cerca de 1.000 palavras, a fim de estimar a cobertura lexical dos recursos construídos. A avaliação foi feita manualmente com base nas palavras reconhecidas e nas não identificadas pelos recursos lexicais aqui construídos.

Palavras-chave

Persa (Farsi), verbo, dicionário eletrónico, morfologia flexional, autómatos e transdutores de estados finitos

Extended Abstract

This project is to be considered as a first step towards the construction of large-scale lexical resources for Persian, and the development of a Persian module to be distributed with the Unitex linguistic development platform.

The Unitex 3.0 has resources for English, Finnish, French (France, Quebec), Georgian (Ancient), German, Greek (Modern and Ancient), Italian, Korean, Norwegian, Polish, Portuguese (Brazil, Portugal), Russian, Serbian (Cyrillic and Latin alphabets), Spanish and Thai.

The specific goal of this dissertation is to build a morphologic, machine-readable dictionary of Persian verbs, using a dictionary of lemmas and a set of morphologic finite-state transducers (FST) to generate all the inflected forms associated to each lemma, and encode them with all the relevant morphosyntactic information (tense, person-number, etc.) because, there are more than 110 million Persian speakers in the world, but the lexical resources for Natural Language Processing (NLP) of the Persian language are still too scarce.

Though Persian (or Fārsi) uses an adapted form of the Arabic script, it is an Indo-European language and its verbal inflection is based on stems and affixes.

This task is complicated in Persian verbal morphology by the fact that each verb has two stems (past and present), different inflection paradigms are used for written (formal) and oral (informal) language uses, and several compound tenses can be formed through combining prefixes and suffixes with base inflected forms, the same tense being able to constitute one, two or more different tokens (separate written forms).

A small dictionary of lemmas (145) and their respective written (292) and spoken (127) stems was built, each stem was provided with the appropriate inflection conventional code, which correspond to an inflection paradigm.

Since the inflected forms of a given verb are produced from its stems, the lemma proper was encoded as a supplementary feature in the verbal entry, which allows for the retrieval of any inflected verb form from texts using that feature.

The list of lemmas was compiled based on frequency data from a large Persian corpus, the TEP (Tehran English-Persian Parallel Corpus), containing around 4.5 million words, and selecting the most frequent verb forms.

For each inflection paradigm, a morphologic FST was built to generate the corresponding inflected forms.

At its current state, the dictionary of inflected forms contains 1,536 entries.

In Persian, in average, each verb lemma yields 28 simple inflected form, 14 simple 'written' inflected forms and 14 'spoken' inflected forms, and relatively, for each stem has 7 inflected forms.

The recognition of compound tenses is carried out by a set of 22 FST using the system morphological mode and the previous lexical annotation of simple verb forms. These FSTs allowed for the retrieval of 3,953 compound verb tenses from the corpus.

Two Backus-Naur form (BNF), one for simple and another for compound inflected forms' entries were built in the form of finite-state graphs and used to automatically validate the format of the lexical entries.

For the evaluation of this language resource, a sample text, retrieved from the www.persian.euronews.com, and containing around 1000 words was used to assess the lexical coverage of the simple words' dictionary and the compound tenses' lexical graphs.

The evaluation was done manually, based on the recognized words and those that were not identified by lexical resources built here.

Only one simple verb form, from the 43 verb forms found (107 occurrences), was not in the dictionary by a lapse in the construction of the respective graph flexion, which has been corrected.

Regarding compound verb forms, the lexical graphs produced 31 lexical entries corresponding to 20 different forms, occurring 46 times in the text. Only one form, involving the contraction of a negation morpheme with a verb (a single instance) was not identified.

Resumo Alargado

Há mais de 110 milhões de falantes da língua Persa (ou Farsi) no mundo, mas os recursos lexicais para Processamento de Língua Natural (NLP) disponíveis para esta língua ainda são muito escassos. Embora o Persa use uma forma adaptada do alfabeto árabe, é uma língua indo europeia e sua flexão verbal é baseada em morfemas radicais e afixos.

Este projeto deve ser considerado como um primeiro passo para a construção de recursos lexicais em larga escala para o Persa, e o desenvolvimento de um módulo do Persa a ser distribuído com a plataforma de desenvolvimento linguístico UNITEX.

O objetivo específico deste trabalho é construir um dicionário eletrónico, legível por máquinas, de verbos do persa, usando um dicionário de lemas e um conjunto de transdutores de estados finitos (FST) morfológicos para gerar todas as formas flexionadas associadas a cada lema, codificando-as com toda a informação morfossintática relevante (pessoa-número, tempo-modo, etc.).

Esta tarefa é complicada na morfologia verbal persa pelo facto de cada verbo tem dois estemas (ing. *stem*), um para os tempos do passado e outro para os do presente; são também empregues diferentes paradigmas de flexão na linguagem escrita (mais formal) e na oralidade (informal); e, finalmente, pelo facto de vários tempos compostos poderem ser formados através da combinação de prefixos e sufixos com formas flexionadas de base, ao mesmo tempo que podem constituir-se num único, em dois ou até mais diferentes palavras gráficas (*tokens*).

Um pequeno dicionário de lemas (145) e seus respectivos estemas da forma 'escrita' (292) e 'falada' (127) foi, então, construído, e cada estema acompanhado pelo código de flexão convencional de flexão adequado, e que corresponde a um paradigma de inflexão.

Uma vez que a flexão é produzida a partir dos estemas, o lema do verbo é tratado como um traço suplementar nas entradas do dicionário, o que permite pesquisar as ocorrências de um dado verbo num texto a partir do respetivo lema.

. A lista dos lemas foi compilada com base em dados de frequência de um grande corpus persa, o TEP (Tehran English – Persian Parallel Corpus), contendo cerca de 4,5 milhões de palavras, e selecionando as formas verbais mais frequentes.

No seu estado atual, o dicionário de formas flexionadas contém 1.536 entradas. Em Persa, em média, cada lema verbal produz 28 formas simples: 14 formas simples 'escritas' flexionadas e 14 formas 'faladas', dado que cada stema.

O reconhecimento de tempos compostos é realizado por um conjunto de 22 FSTs utilizando o modo morfológico do sistema e a prévia anotação lexical de formas verbais simples. Estes grafos permitiram reconhecer 3.953 formas compostas no corpus utilizado.

Duas formas Backus-Naur (BNF), uma para as entradas das palavras simples e outra para as entradas de formas compostas flexionadas, foram construídas sob a forma de gráficos de estados finitos e utilizadas para validar automaticamente o formato das entradas lexicais.

Para a avaliação deste recurso linguístico, utilizou-se um texto obtido a partir da www.persian.euronews.com e contendo cerca de 1.000 palavras, a fim de estimar a cobertura lexical dos recursos construídos.

A avaliação foi feita manualmente com base nas palavras reconhecidas e nas não identificadas pelos recursos lexicais aqui construídos.

Apenas uma forma simples de um verbo, de entre as 43 encontradas (107 ocorrências), não estava registada no dicionário, por lapso na construção do respetivo grafo de flexão, o que foi corrigido.

Quanto às formas verbais compostas, os grafos lexicais produziram 31 entradas lexicais, correspondendo a diferentes flexões de 20 formas, que ocorrem 46 vezes no texto.

Apenas uma forma, correspondendo à contração de um morfema de negação com um verbo (uma única ocorrência) não foi identificada.

TABLE OF CONTENTS

1. Introduction	1
1.1. Objectives.....	2
2. Methodology.....	3
2.1. Synopsis of Persian verb morphology.....	3
2.2. A lexicon of Persian verbs	5
2.3. A Persian verbs dictionary of inflected forms.....	7
2.3.1. BNFs	8
2.3.2. Inflection FSTs for Written simple words.....	11
2.3.3. Inflection FSTs for Spoken simple words	13
2.3.4. Compound inflection FSTs	17
2.4. Generating the dictionary of inflected simple forms.....	20
2.5. Generating the inflected forms.....	22
2.6. Apply the DELAF to the corpus	23
2.7. Apply the compound tenses FST to the text.....	24
3. Evaluation	25
3.1. Evaluation of the formal correction of the dictionary	25
3.2. Evaluation of the lexical coverage of the dictionary.....	25
3.2.1. Lexical coverage of simple forms' dictionary.....	26
3.2.2. Lexical coverage of compound forms' dictionary	26
4. Conclusion.....	28
References	30
Appendices.....	31
A. DELAS (dictionary of Persian verb stems and associated lemmas)	31
B. Sample of DELAF (dictionary of inflected simple forms)	35
C. Sample of Compound forms (found in a text)	36
D. Simple BNF and Compound BNF.....	37
E. Sample of simple verbs FST.....	38
F. Sample of Compound verbs FST	39
G. Evaluation sample text.....	40

LIST OF FIGURES

Figure 1. BNF1 for simple tenses	8
Figure 2. BNF2 for compound tenses.....	10
Figure 3. V001: written simple present dictionary graph	11
Figure 4. V002: written simple past dictionary graph.....	12
Figure 5. V003: است ast (be) dictionary graph	12
Figure 6. V004: هست hast (be) dictionary graph.....	13
Figure 7. V005: spoken simple present dictionary graph without root.....	14
Figure 8. V008: spoken simple present dictionary graph with root	14
Figure 9. V007: spoken simple present dictionary graph with changes in HS3sm	15
Figure 10. V009: spoken simple past dictionary graph.....	16
Figure 11. V006: spoken simple past dictionary graph with changes in GS2pm and GS3pm.....	16
Figure 12. Compound-tenses dictionary graph	17
Figure 14. GC sub-graph: Past continuous (GC).....	18
Figure 15. HU sub-graph; Present Subjunctive dictionary graph.....	19
Figure 16. Extract of Persian DELAS	20
Figure 17. Extract of Persian DELAF	22
Figure 18. Word List.....	23
Figure 19. First entries of the compound tenses dictionary	24

LIST OF TABLES

Table 1. Lexicon table of Verbs	6
Table 2. The distribution of the verb stems by inflection graphs	21

1. Introduction

There are approximately 110 million Persian (natively known as فارسی **Farsi** [fɒ:r'si:] or پارسی **Parsi**) speakers worldwide¹, with the language property official status in Iran, Afghanistan and Tajikistan. For centuries, Persian has likewise been a prestigious cultural language in Central Asia, South Asia, and Western Asia. Persian is used as a liturgical language of Islam not only in Iran, Afghanistan, and Tajikistan, but also in Pakistan and North India.

In spite of its cultural and demographic importance, studies in Natural Language Processing (NLP) on Persian are scarce. One of the main problems is the fact that lexical resources for NLP are almost non-existent in Persian, as well as, no literature review existed for this theme.

We intend to start the building of a new module for the Unitex² linguistic development platform (Paumier 2003, 2013) by describing the inflection of Persian verbs.

This project is to be measured as a first step towards the construction of large-scale lexical resources for Persian, and the development of a Persian module to be distributed with the Unitex linguistic development platform.

The project consists in formalizing the main inflectional paradigms of Persian verbs and produce automatically all the inflected forms, simple and compound, both formal and informal forms³, associated with them from a dictionary of lemmas.

¹ http://en.wikipedia.org/wiki/Persian_language [URL <access in 10-05-2013>]

² <http://www-igm.univ-mlv.fr/~unitex/>

³In Persian, verbs can be used both in formal and informal speech and these informal uses occur most of the times in speaking but they can appear in writing also; formal and informal verb forms are morphologically different so they have to be taken into account in inflectional dictionary. In following pages, we call these verbs as a *spoken* (informal) and *written* (formal) forms. For example, آمد *amad* is the written stem of the verb "come" but the spoken stem is اومد *omad*.

Unitex is a collection of programs developed for the analysis of texts in natural language by using linguistic resources and tools. These resources consist of electronic dictionaries, local grammars and the lexicon-grammar tables with lexicon-grammatical information about the predicative element's properties of a language (M. Gross 1996).

These languages already exist in Unitex: English, Finnish, French (France, Quebec), Georgian (Ancient), German, Greek (Modern and Ancient), Italian, Korean, Norwegian, Polish, Portuguese (Brazil, Portugal), Russian, Serbian (Cyrillic and Latin alphabets), Spanish and Thai.

In such electronic dictionary of inflected forms, each verb form is linked to the corresponding lemma, its part of speech (POS) and all grammatically relevant information pertaining to that word.

1.1. Objectives

For achieving the aim of this project, we must perform these following steps:

1. Build a sufficiently large lexicon of Persian verbs, in parliamentary procedure to reach a reasonable lexical coverage;
2. Produce a Persian verbs dictionary of lemmas, associated with their corresponding inflectional paradigms;
3. Build inflectional graphs to generate automatically the both simple and compound forms of Persian verbs, in order to produce a dictionary of inflected forms, aiming at a highly granular and systematic description of all the relevant grammatical values pertaining to those forms.
4. Apply the BNFs to validate automatically the inflect words' dictionary that retrieved from the corpus to evaluate the formal correction of the result.
5. Apply the Persian dictionary to a sample of text in order to assess the lexical coverage and grammatical adequacy of the dictionary.

2. Methodology

2.1. Synopsis of Persian verb morphology

Persian has a complex verbal tense system⁴ (Tabatabayi, A. (2010))⁵ (its Instructional Technology Services. (2007))⁶. Besides having specific endings for some tenses, for instance:

- Simple Past (GS⁷), رفتم⁸ *raftam*, (lit. “went I”), means “I went”
- Simple Present (HS) روم, *ravam*, (lit. “go I”), means “I go”;

Certain tenses are produced by adding prefixes to some inflected forms, for example:

- Past Continuous (GC), می رفتم (میرفتم) *miraftam [mi+raftam]*, (lit. “ing=duration went I”), means “I was going” which is formed by the prefix *mi-* (می) and the simple past forms followed by the suffix *-am* (آم) ;
- Present Continuous (HC) می روم (میروم), *miravam [mi ravam]*, and (lit. “ing=duration go I”), means “I am going” that is also shaped by the prefix *mi-* (می) and the simple present forms followed by the suffix *-am* (آم) . These two tenses can be composed with the prefixes both joined and parted;

⁴ http://en.wikipedia.org/wiki/Persian_verbs [URL <access in 15-06-2013>]

⁵ <http://dlib.ical.ir/site/catalogue/789489> [URL <access in 15-06-2013>]

⁶ http://sites.la.utexas.edu/persian_online_resources/verbs/ [URL <access in 01-07-2013>]

⁷ The codes here used for tenses are conventional, but they were derived from the Persian terminology, for easy mnemonic; **GS**= Gozashteye Sade (گذشته ساده) Simple Past, **HS**= Hale Sade (حال ساده) Simple Present, **GC**= Gozashteye estemrari (گذشته استمراری) Past Continuous, **HC**= Hale estemrari (حال استمراری) Present Continuous, **HP**= Hale kamel (حال کامل) Present Perfect, **HU**= Hale eltezami (حال التزامی) Present Subjunctive, **I**= amri (امری) Imperative, **HF**= Hale Kamel estemrari (حال کامل استمراری) Present Perfect Continuous, **GP**= Gozashteye kamel (گذشته کامل) Past Perfect, **GF**= Gozashteye Kamel estemrari (گذشته کامل استمراری) Past Perfect Continuous, **HK**= hale kamele eltezami (حال کامل التزامی) Perfect Subjunctive), **AS**= Ayande (آینده) Simple Future.

⁸ All example provided in a first person singular.

- Present Subjunctive (HU) بروم *berravam*, (lit. “be=subjunctive particle go I”), means, “(that) I go” that is formed by the prefix *be-* (بـ) and the simple present forms followed by the suffix *-am* (آم).
- Imperative (I) برو *boro*, (lit. “b=imperative particle go you”), means “(you) GO!” which is made by the prefix *b-* (بـ) followed by the root of the present form *ro* (رو);

In addition, Persian has one tense, which is created by appending suffixes to the simple past form of the verb like as:

- Present Perfect (HP), ام رفته *rafteh am*, (lit. “gone I”), means “I have gone” which is made by the simple past forms followed by the suffix *-ham* (آم ه);

Furthermore, compound tenses are created by adding a prefix and a suffix to some inflected forms and they can be written as two tokens or three tokens compound, such every bit:

- Present Perfect Continuous (HF) ام می رفته *mirafteh am* (میرفته ام), *mirafteh am* [*mi+rafteh+am*], and (lit. “ing=duration gone I”), means “I had been going” is formed by the prefix *mi-* (می) and the simple past forms followed by the suffix *-ham* (آم ه).

There are also some compound tenses made by adding suffix or prefix to some inflected forms, but in this case, the suffix or prefix is not meaningless words, it is a verb and also for changing, the person of the verb the suffix or prefix verb changed not a basic verb, such as:

- Past Perfect (GP) بودم رفته *rafteh bodam* (lit. “gone been I”), means “I had gone” is formed by the simple past third singular person followed by *-h* (ه) and inflected form of the simple past forms of *bod* (بود) *been*;

- Past Perfect (GF) *رفته بوده ام*, *rafteh bode am* [*rafteh+bode+am*], (lit. “gone had been I”), means “I had gone” is formed by the simple past third singular person followed by *-h* (ه) and the simple past forms of *bod* (بود) *been* followed by the suffix *-ham* (ه ام);
- Present Perfect Subjunctive (HK) *رفته باشم*, *rafteh basham* (lit. “gone be I”), means “I have gone” is formed by the simple past third singular person followed by *-h* (ه) and inflected forms of the simple present forms of *bash* (باش) *be*;
- Simple Future (AS) *خواهم رفت*, *khaham raft* (lit. “Want I go”), means “I will go” is formed by inflected forms of the simple present forms of *Khah* (خواه) *want* followed by the simple past third singular person form of the verb.

2.2. A lexicon of Persian verbs

In order to select the most frequent verbs for the lexicon, we used the TEP corpus: Tehran English-Persian Parallel Corpus (Pilevar *et al.* (n.d.))⁹ obtained from the Natural Language and Text Processing Laboratory of Tehran University¹⁰. The corpus consists of 4 million tokens and, after being tokenized with Unitex 3.0, it features 556,234 sentences, 15,166,987 (64,492 different) tokens, 4,485,147 (64,365 different) simple word forms and 3,239,250 (10 different) digits.

Using the frequency list of the corpus’ words ranked by decreasing frequency and automatically listed by the Unitex, we estimated the cumulative frequency of its word forms.

We then annotated the verb forms, indicating their lemma, starting from the most frequent words until we attained a rank corresponding to a cumulative percentage of 90% of the total corpus word forms (around 3 million and a half) corresponding to 1,466 different verb forms.

⁹ <http://ece.ut.ac.ir/nlp/resources.htm>

¹⁰ <http://ece.ut.ac.ir/nlp/>

Count	Token	Cum.count	Cum. %	Lemma	Stem	Transliteration	Translation	V. Class	DELAS entry (Stem)
16877	بود	16877	0.015707	بودن	بود	bod	Been	V002	V002, بود
12935	باشه	29812	0.027746	بودن	باش	bash	Be	V001	V008, باش
11564	شده	41376	0.038509	شدن	شد	shod	became	V002	V002, شد
10835	کنم	52211	0.048593	کردن	کن	kon	Do	V001	V001, کن
10295	داره	73057	0.067994	داشتن	دار	dar	have	V001	V008, دار
9428	کن	82485	0.076769	کردن	کن	kon	do	V001	V001, کن
9344	دارم	91829	0.085466	داشتن	دار	dar	have	V001	V001, دار
8614	است	100443	0.093483	بودن	است	ast	Is	V003	V003, است
7104	هست	107547	0.100094	بودن	هست	hast	Is	V004	V004, هست
6932	داري	114479	0.106546	داشتن	دار	dar	have	V001	V001, دار
6756	میشه	121235	0.112834	شدن	شو	sho	become	V001	V008, شو

Table 1. Lexicon table of Verbs

Table 1 indicates a fragment of the lexicon of verbs, built from the list of simple row and their frequencies, where just the verb classes are presented, in decreasing frequency in decreasing frequency.

The Count column shows the frequency of the token in the corpus and the Cumulative count column presents the total number of verb forms up to that word.

The Cumulative percentage column indicates the percentage of the corpus that the previous Cumulative count represents.

For each token, the corresponding lemma was added. Since, in Persian, the simple tenses are derived from two distinct stems, each stem constitutes a different entree, but the same lemma is provided for both stems, so that the two sets of inflected forms will be linked up to the same lemma.

For example: رفتن *raftan* “to go” has two stems; رو *ro* “go” is the stem of the present form and رفت *raft* “gone” is the stem of past form of the same verb. In Persian, lemma builds by adding “ان *an*” to the end of past stem, for example: رفت *raft* + ان *an* = رفتن *raftan*.

Also for more consideration about each verb stem and the lemma, we use two websites (Redirected from Persian in Texas)¹¹ and (فعل ساده *feelee sade* “Simple Verb”¹²) that these websites are contain of Persian verb stems and lemmas.

2.3. A Persian verbs dictionary of inflected forms

For building a dictionary of lemmas and their respective written and spoken stems, for each stem, we must provide the appropriate inflection conventional code, which correspond to an inflection paradigm.

Since the inflected forms of a given verb are produced from its stems, the lemma proper was encoded as a supplementary feature in the verbal entry, which allows for the retrieval of any inflected verb form from texts using that feature.

¹¹ <http://persian.nmelrc.org/pvc/simpleverbs.php>

¹² http://fa.wikipedia.org/wiki/%D9%81%D8%B9%D9%84_%D8%B3%D8%A7%D8%AF%D9%87

2.3.1. BNFs

In order to provide formal guidelines for the description of the verb tense inflection¹³, we have built a BNF (Backus Normal Form or Backus–Naur Form)¹⁴ with all the relevant grammatical information pertaining to the verb inflection. In computer science, a BNF is an annotation technique for context-free grammars, often used to identify the syntax of languages used in computing, often used to describe the syntax of languages used in computing. Two Finite-State Automata (FSA), shown in **Figure 1** and **Figure 2**, were built to be used as BNFs for the validation of the simple and compound inflected words' dictionaries.

The goal is to use these FSA in order to validate automatically the inflect words' dictionary entries.

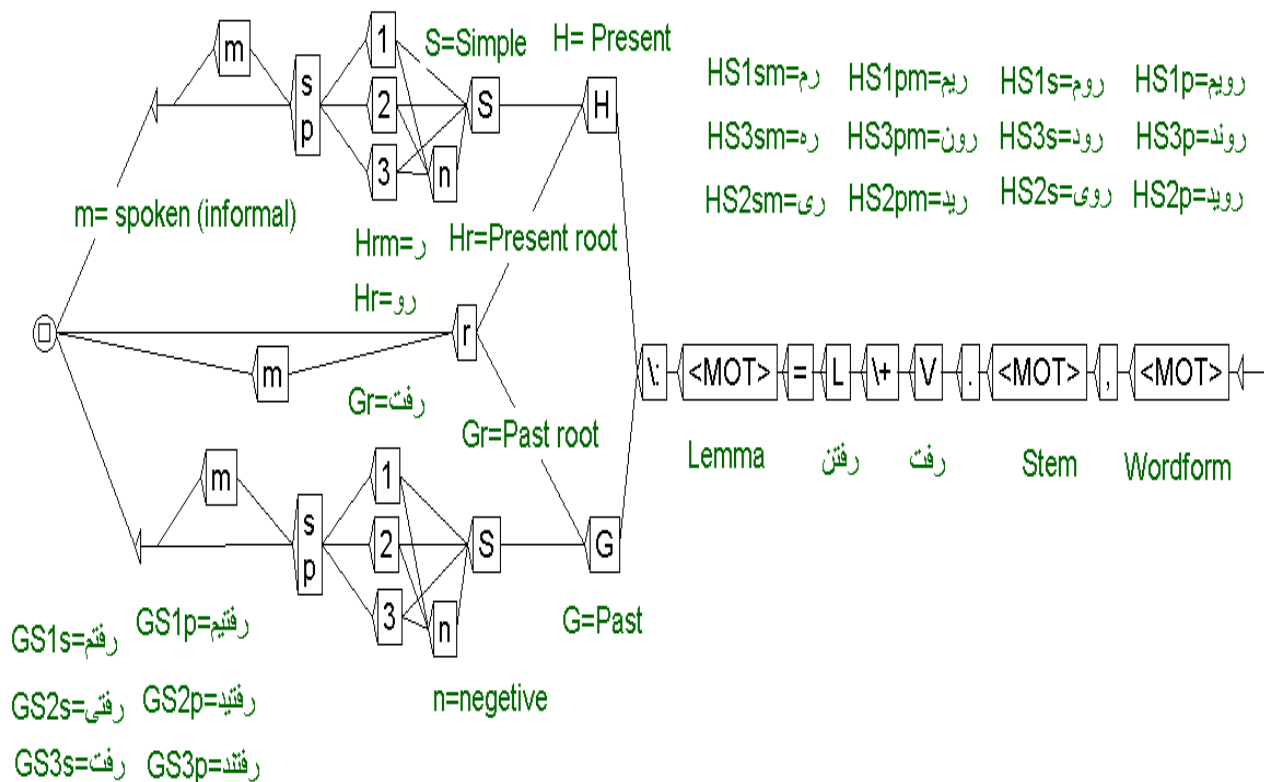


Figure 1. BNF1 for simple tenses

¹³ http://en.wikipedia.org/wiki/Persian_verbs

¹⁴ http://en.wikipedia.org/wiki/Backus%E2%80%93Naur_Form

The graph (**Figure 1**) describes the format of an inflected word entry in the DELAF (dictionary of inflected form) and it reads (right to left¹⁵) as follows: the first box with the in-built mask <MOT> stands for any simple word (one token); it is followed by a comma ‘,’ and the second word stands for the stem; then a dot ‘.’ separates it from the part-of-speech and here, only “V” for verb, followed by ‘+’ and “L” is a code for lemma, “=” shows that the “L” is equal to the corresponding Lemma that shows as a third word (MOT); in this stage, the verb entry has two major paths issuing from that point onward: one for “H”, representing the present forms, and another for “G”, which represents the past forms; a third path goes to “r” for the roots of both the present and the past forms. As this is a simple word entries’ BNF, the past and the present forms have only simple tenses; thus, each main path is followed by the “S” for simple tense, the code “n” is used to show the negation, the “1”, “2”, and ”3” symbols for 1st, 2nd, 3rd person respectively, “s” and “p” for singular and plural forms which, “m” stands for spoken (informal) verb forms and also the line without “m” shows the written (formal), for example: *HSIs* is written simple present first singular person and *HSIsm* is spoken simple present first singular person. The same is done for the past forms.

The graph in the following figure (**Figure 2**) is read in the similar way as the previous one, but it describes the entries of compound tenses. The first box shows a loop for the compound form entry (more than one token). For the present tenses “H”, we find the codes¹⁶ for the Continuous “C”, Subjunctive “U”, Perfect Continuous “F”, Perfect “P” and Perfect Subjunctive “K” tenses; for the past tense “G”, we find the Continuous “C”, Perfect continuous “F” and Perfect “P”. For the Imperative mode the code “I” is used, but only the second person value has been considered, both

¹⁵ Normally, graphs in Unitex are from left to right but the direction here is opposite because writing order in Persian is from right to left.

¹⁶ We do not use the same letters (codes) as in Unitex manual, because it does not have the all letters which we need, also we wanted that the present (H) and the past (G) be divided to be more understandable.

for singular and plural. For the Future tense “A”, there is only the Simple “S”, which in fact is a compound form (e.g. *خواهم رفت*, *khaham raft* (lit. “Want I go”) means, “I will go”), the “simple” being the traditional grammatical term used to describe it.

Examples of all inflections are provided for the verb *رفتن* “to go”.

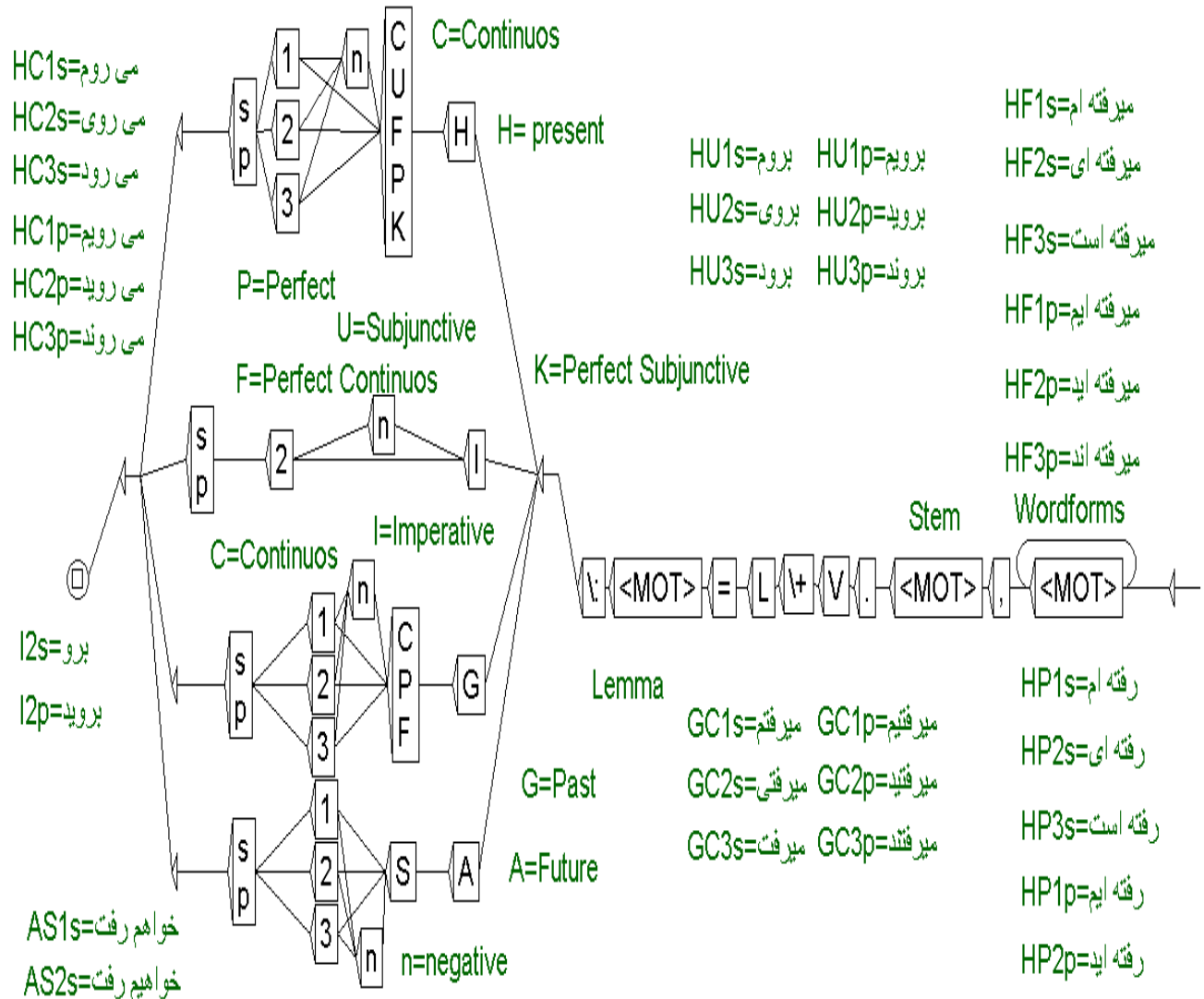


Figure 2. BNF2 for compound tenses

2.3.2. Inflection FSTs for Written simple words

In order to generate the inflected forms associated to each lemma, we built four graphs for written simple forms of verbs. These graphs are shown in **Figure 3** and **Figure 4**.

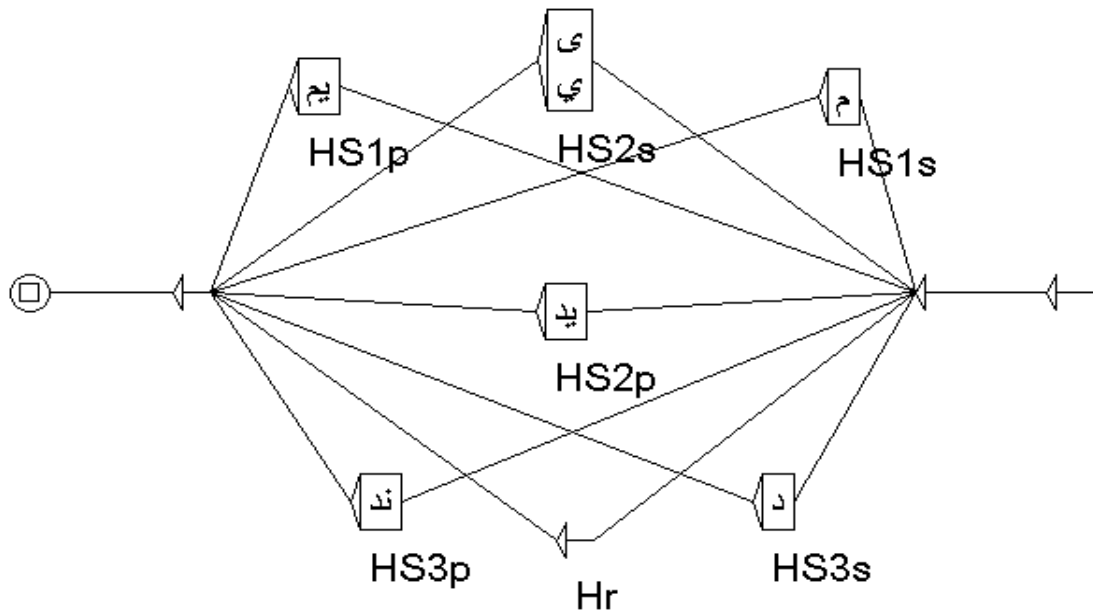


Figure 3. V001: written simple present dictionary graph

This graph describes changes the stem of the word undergoes in order to generate all inflected forms of written simple present. Starting with the box in the upper-right side, we see that the simple present first person singular (HS1s) is produced by adding “*m*” to the stem, for example: for the verb رفتن *raftan* “to go” whose stem of the present is *رو rav*¹⁷, this will yield the form *روم ravam*.

¹⁷ When the present stem of the verb ends up with *-av* as in *رو rav* “go”, it changes to *-o* as in the present root (Hr) *رو ro*. In the second person singular of imperative form prefix *بی be-* also becomes *بُو bo-*.

The same is done for the remaining inflections in **Figure 4**; the graph shows the inflection of the written simple past.

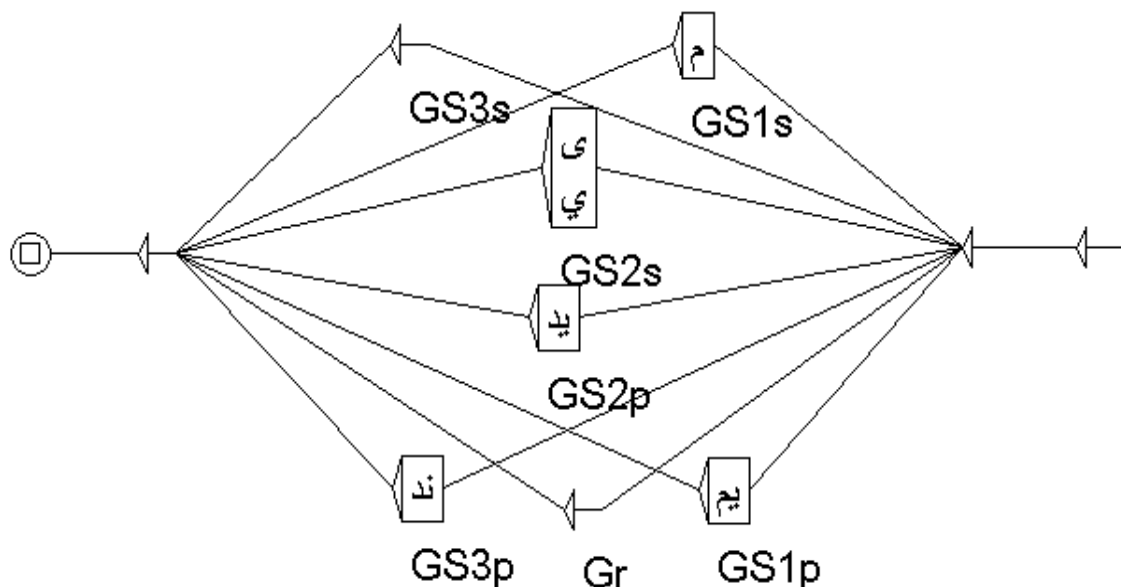


Figure 4. V002: written simple past dictionary graph

In Persian, we have two verbs *است / ast* “to be” and *هست / hast* “to be” that derived from the verb of *بودن / bodan* “to be” which do not follow the same structure as simple present. Because of this, for *ast* “be” and *هست / hast* “be”, a specific set of graphs must be built. These are shown in **Figure 5** and **Figure 6**.

In the case of *است / ast* “be”, only the present simple third singular person (HS3s) and the present root (Hr) exist.

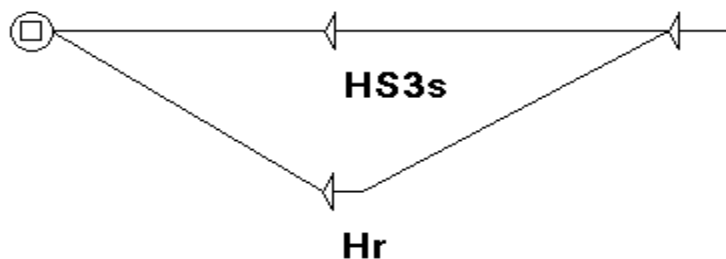


Figure 5. V003: است ast (be) dictionary graph

In the case of *هست hast* “be”, only the simple present (HS) tenses (all forms), and the present root (Hr) exist, the difference between this verb inflections from the other simple present verb inflections is that, this verb followed the simple past inflections rule.

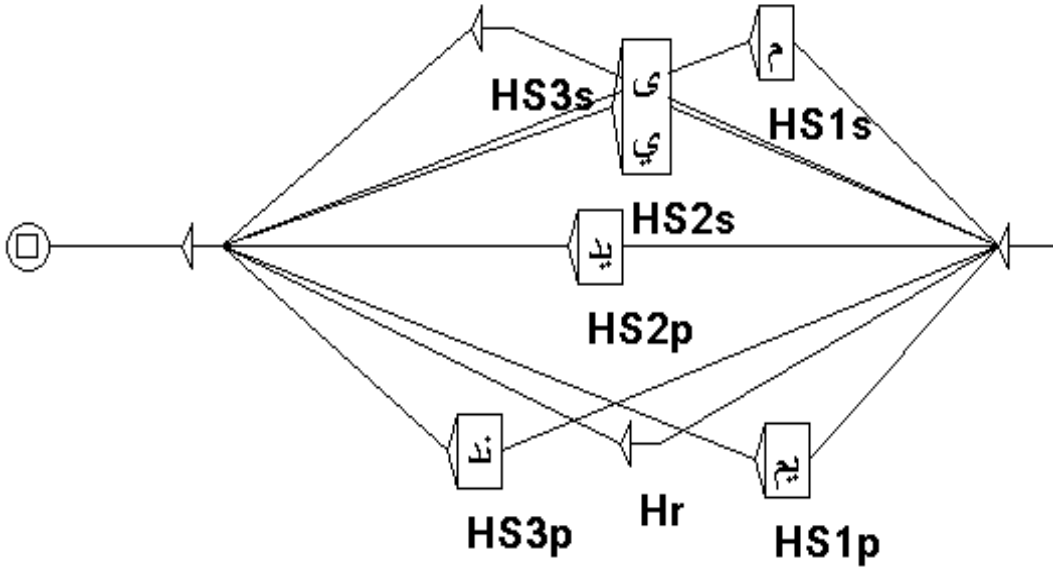


Figure 6. V004: *هست hast* (be) dictionary graph

2.3.3. Inflection FSTs for Spoken simple words

As noted before, Persian uses different inflections for written (formal) and spoken (informal) forms.

Nevertheless, these forms also appear in written texts. So, they also have to be taken into account. In order to generate the inflected forms associated to each lemma of the spoken inflections, we built five graphs. These graphs are shown in **Figure 7**, **Figure 8**, **Figure 9**, **Figure 10** and **Figure 11**, which explain the differences between these five graphs below. The “spoken” forms are signaled by the code ‘m’. The term *محاوره ای mohavereyi* “spoken”, used in Persian grammar to denote the spoken form.

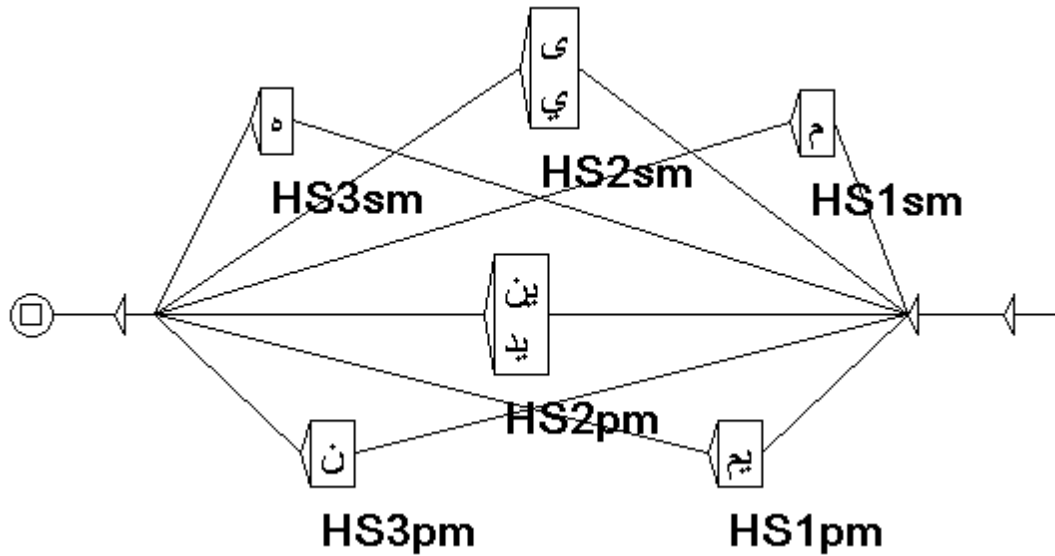


Figure 7. V005: spoken simple present dictionary graph without root

The graph in the **Figure 7** is read in the similar way as the written simple present graph, but only suffixes in plural forms are different and also it does not have a root.

We made another graph for a spoken simple present, which are presented in the graph of **Figure 8**; simply, the difference of this graph from previous one is that this graph has a root form (Hrm).

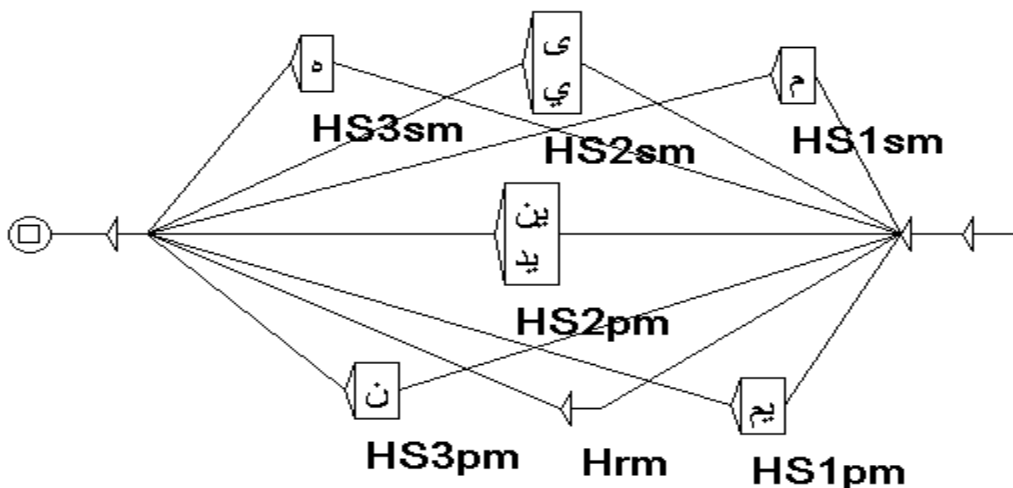


Figure 8. V008: spoken simple present dictionary graph with root

In the other cases, some verbs get *d* instead of *h* as a suffix in the third singular person (HS3sm), for example: the verb *khastan* “want” with the spoken stem of *kha* “want” in a third singular person from gets *d* (*kha* + *d* = *khad* *خواد*).

This graph is shown in **Figure 9**.

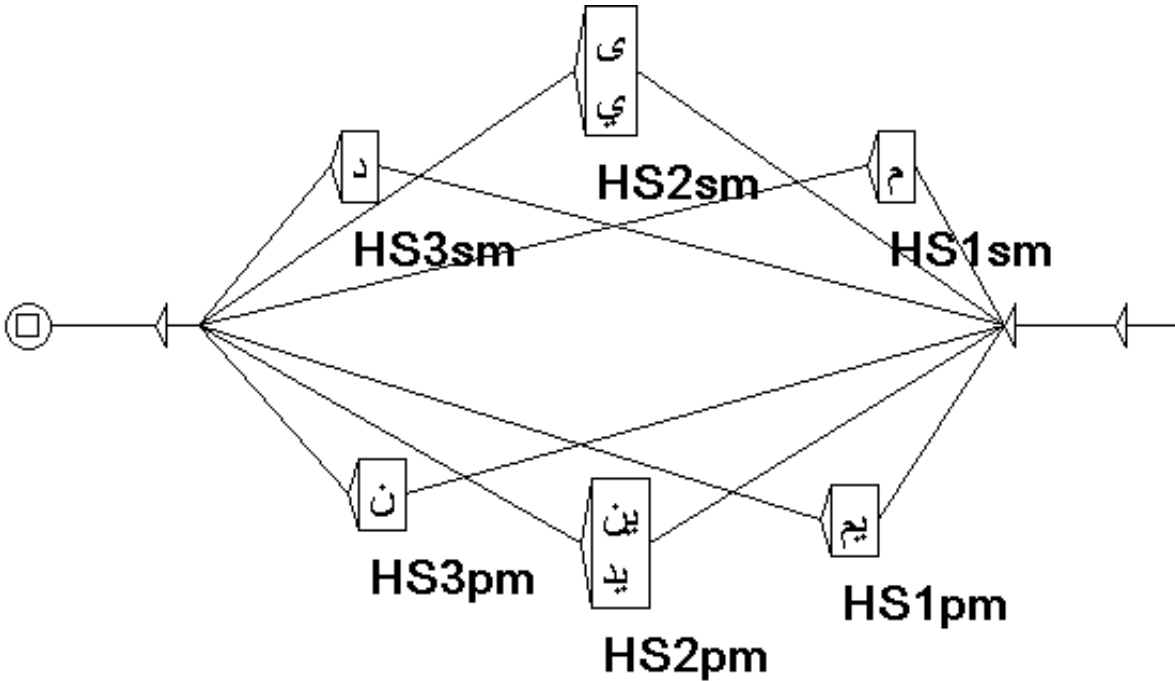


Figure 9. V007: spoken simple present dictionary graph with changes in HS3sm

The graph in the **Figure 10** is read in the similar way as the written simple past graph, but only suffixes in plural forms are changed, this graph made and used only for a spoken (informal) stem that never use these stems in formal writing.

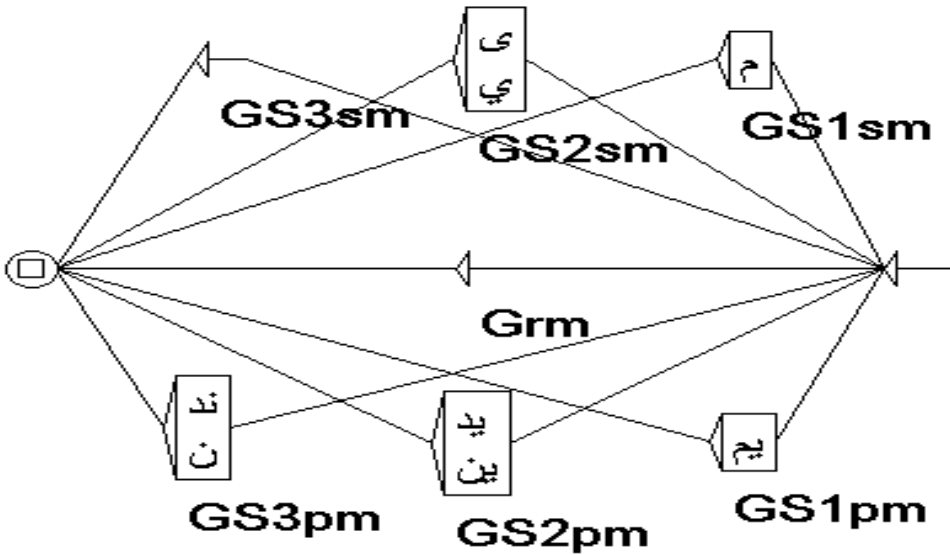


Figure 10. V009: spoken simple past dictionary graph

We added one more graph in **Figure 11** for spoken simple past, sometimes the past written (formal) stem of the verb by getting a spoken (informal) suffixes become an informal verb and this case only happened in second and plural person of simple past. For example: the verb پختن *pokhtan* “cook” with stem of پخت *pokht* “cook” in written (formal) third plural person get کند *and* (پخت *pokht* + کند *and* = پختند *pokhtand*) but if we add ن *an* (informal suffix) instead of کند *and*, it has become an informal verb (پخت *pokht* + ن *an* = پختن *pokhtan*).

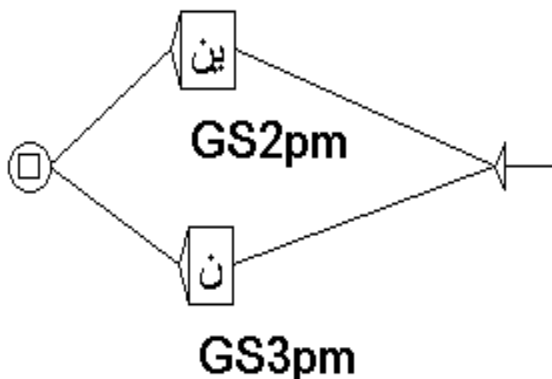


Figure 11. V006: spoken simple past dictionary graph with changes in GS2pm and GS3pm

2.3.4. Compound inflection FSTs

Persian adopts a complex system of compound tenses to express other mode, tense, and aspect variation (See 2.1).

We built several graphs for compound tenses as sub-graphs¹⁸ in one main graph, which intends to capture all correspond inflected forms. This is shown in **Figure 12**.

The sub-graph on the left consist of the compound tenses without negation, while the sub-graphs on the right describe the negative compound tenses. This is indicated in the graphs names by an “n”.

For example, **HC** stands for the Present Continuous, while, **HCn** designates the negative Present Continuous.

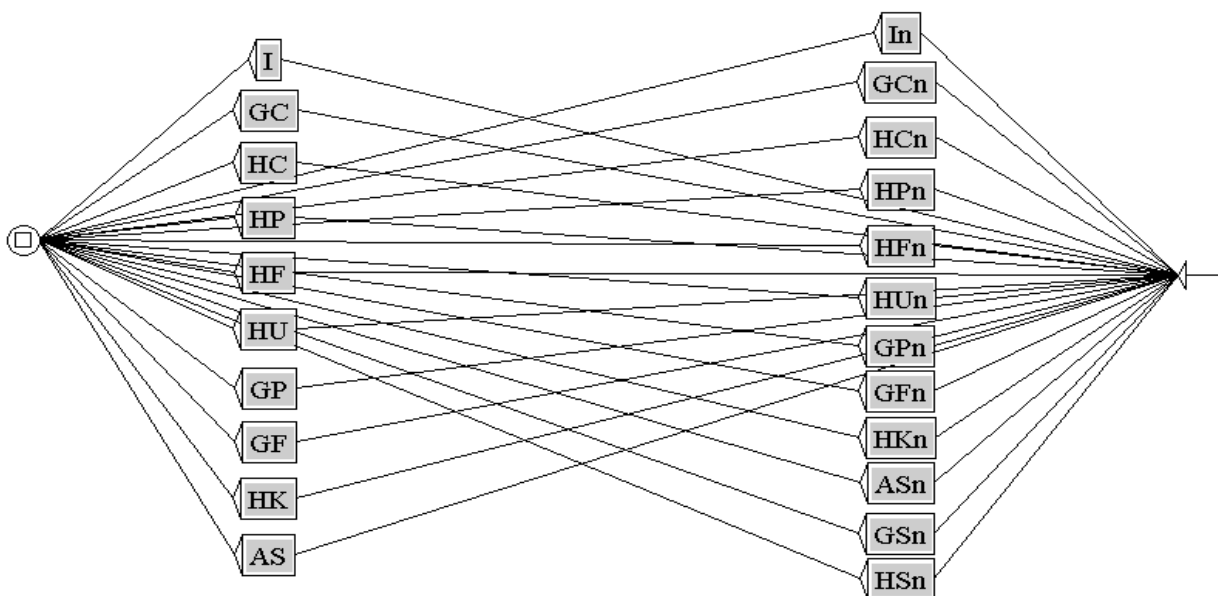


Figure 12. Compound-tenses dictionary graph

To illustrate a sub-graph, consider **Figure 13**.

¹⁸ Each sub-graph is one graph, that we put it in one main graph by using a “:” and the name of the graph, for example: we write “:HC” in one box and Unitex automatically find the graph, you can open a sub-graph by clicking on the grey line while pressing the Alt key.

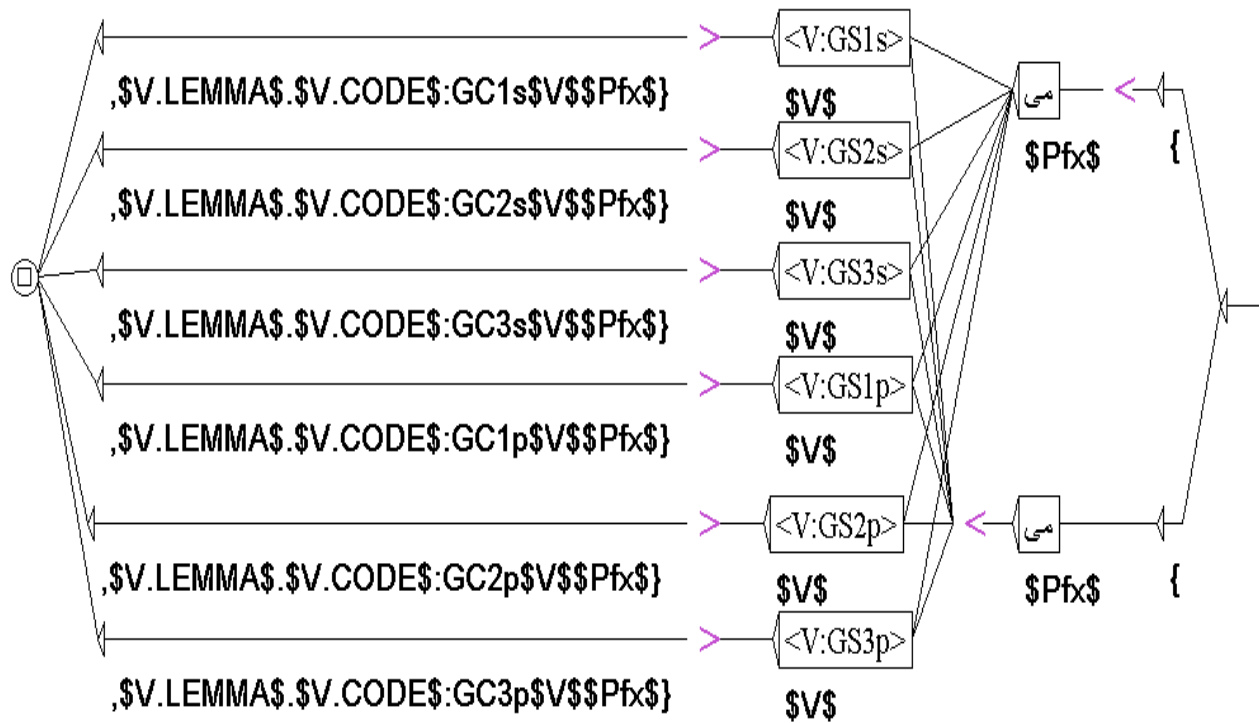


Figure 13. GC sub-graph: Past continuous (GC)

This graph describes the past continuous compound forms and it has two main paths that correspond to the possibility of the prefix *می* *mi-* to be attached (upper path) or split up (lower path) from the base word. In the first path, the simple form that is to be recognized is marked between > < to indicate that the graph is working in the morphological mode (Paumier 2013 page 67 ff.; 129 ff.). The string *می* *mi-* is associated to the variable *\$Pfx\$* because it is a prefix, which is then followed by the form that has already been tagged as <V: GS1s>, that is a verb in the past simple first person singular. This form is associated to variable *\$V\$*. In the output, the system produces the compound entry which consist of the content of the variables *\$Pfx\$* and *\$V\$*, the lemma of *\$V\$*, the codes of *\$V\$* and the inflection code for the past continuous first singular person {GC1s}. A similar procedure is carried out for the remaining simple inflected forms. Along

the second path, the only difference is that *می* *mi-* is taken care of as a separate word but the output is the same.

Moreover, we write an output in this order because the graph is applied from right to left, in the correct order, so that the output must be something like: {*\$Pfx\$\$V\$, \$V. LEMMA\$. \$V.CODE\$:GC1s*}.

As an analyzed form, these compound inflected forms are delimited by the meta-symbols “{“ and “}”.

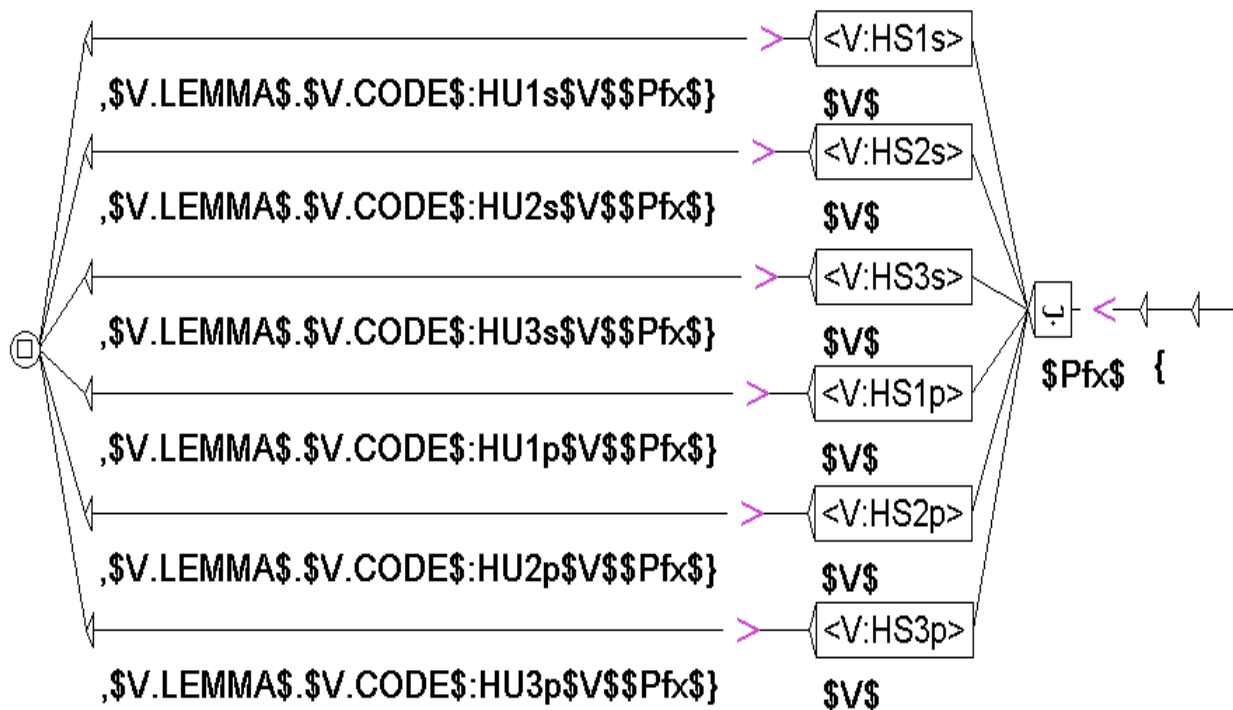


Figure 14. HU sub-graph: Present Subjunctive dictionary graph

Figure 14 is similar to the previous one, only there is just a simple word (one token), regularly derived by adding the prefix *بـ* *be-* to the simple present forms in order to produce the present subjunctive inflected forms.

2.4. Generating the dictionary of inflected simple forms

In order to construct the dictionary of inflected forms, we first have to build a dictionary of lemmas, which has the following structure:

`<stem>,<inflection-code>+L=<LEMMA>`

where `<stem>` represents the form that will be used to generate the inflected forms; and `<inflection-codes>` is composed of the word part-of-speech (POS) in this case “L” for the lemma, and conventional number for indicating the inflectional FST describing that inflection paradigm.

For example, here are some of the entries of the Persian verbs DELAS (Dictionary of Simple Words Lemmas):

رفت، رفتن=L+V002
رو، رفتن=L+V001
آمد، آمدن=L+V002
آی، آمدن=L+V001
شد، شدن=L+V002
شو، شدن=L+V001

Figure 15. Extract of Persian DELAS

The full lexicon is provided in appendix A.

Table 2 shows the distribution of the verb stems by inflection graphs:

Graph	Total stems	Example (3 rd person singular)	Transliteration	Gloss
V001	136	آید	ayad	He/she comes
V002	154	آمد	amad	He/she came
V003	1	است	ast	he/she/it is
V004	1	هست	hast	he/she/it is
V005	7	ده	deh	He/she gives
V006	120	آمده	amadeh	He/she came
V007	3	خواد	khad	He/she wants
V008	19	دونه	doneh	He/she knows
V009	21	اومده	omadeh	He/she came

Table 2. The distribution of the verb stems by inflection graphs

2.5. Generating the inflected forms

By intersecting the DELAS and the inflection FSTs, Unitex is able to generate all the inflected forms associated to each stem of the verbs represented in the dictionary. This constitutes the DELAF (Dictionary of Inflected Forms).

The output is shown in Figure 16:

```
GS1s: رفتن=L+V. رفت ، رفتم
GS2s: رفتن=L+V. رفتی ، رفتی
GS3s: رفتن=L+V. رفت ، رفت
GS1p: رفتن=L+V. رفتیم ، رفتیم
GS2p: رفتن=L+V. رفتید ، رفتید
GS3p: رفتن=L+V. رفتند ، رفتند
Gr: رفتن=L+V. رفت ، رفت
```

Figure 16. Extract of Persian DELAF

A sample of this DELAF is provided in Appendix BB.

Each entry in DELAF has the following structure:

```
<word-form> , <stem> . <POS> + L = <lemma> : <inflection>
```

Where the <word-form> is the inflected form of the verb, <stem> is the stem represented in the DELAS; <POS> is the grammatical category or part of speech (POS) of each word form and <L> is the value of the lemma that is associated with each form.

So far, the dictionary contains 145 lemmas, 292 written verb stems and 127 spoken verb stems that have been described in the DELAS. These allow for the generation of 1,536 inflected forms.¹⁹

¹⁹ We evaluated the formal correctness of the DELAF entries by using the BNF (see 3.1, below)

2.7. Apply the compound tenses FST to the text

When applied compound tenses FST to the corpus, the graph of compound tenses matches 222,698 instances of 3,953 different compound forms. In the merge mode, the output also shows the stem and the lemma of verbs. **Figure 18** illustrates the first entries of the compound tenses: one can see that the first two lines correspond to alternative wordings of present perfect third person singular with and without *است / ast* “is”.

HP3s: کرده ,کرد =V+L.کردن	GP2p: کرده بودید,کرد =V+L.کردن	HP3s: گذاشته ,گذاشت =V+L.گذاشتن
HP2s: کرده ای,کرد =V+L.کردن	GP1p: کرده بودیم,کرد =V+L.کردن	HP3s: گذاشته است,گذاشت =V+L.گذاشتن
HP2p: کرده اید,کرد =V+L.کردن	GP1s: کرده بودم,کرد =V+L.کردن	GP1s: گذاشته بودم,گذاشت =V+L.گذاشتن
HP1p: کرده ایم,کرد =V+L.کردن	GP3p: کرده بودند,کرد =V+L.کردن	GP3p: گذاشته بودند,گذاشت =V+L.گذاشتن
HP3s: کرده است,کرد =V+L.کردن	GF3s: کرده بوده ,کرد =V+L.کردن	GF3s: گذاشته بوده ,گذاشت =V+L.گذاشتن
HP1s: کرده ام,کرد =V+L.کردن	GF3s: کرده بوده است,کرد =V+L.کردن	HP3s: گذاشته ,گذشت =V+L.گذاشتن
HP3p: کرده اند,کرد =V+L.کردن	HP3s: کنده ,کند =V+L.کندن	HP2s: گذاشته ای,گذشت =V+L.گذاشتن
HK2s: کرده باشی,کرد =V+L.کردن	HP3s: کنده است ,کند =V+L.کندن	HP3s: گذاشته است,گذشت =V+L.گذاشتن
HK2p: کرده باشید,کرد =V+L.کردن	GP3s: گذاشته بود,گذاشت =V+L.گذاشتن	HP1s: گذاشته ام,گذشت =V+L.گذاشتن
HK1p: کرده باشیم,کرد =V+L.کردن	GP2s: گذاشته بودی,گذاشت =V+L.گذاشتن	HP3p: گذاشته اند,گذشت =V+L.گذاشتن
HK3s: کرده باشید,کرد =V+L.کردن	HP1s: گذاشته ام,گذاشت =V+L.گذاشتن	HK3s: گذاشته باشد,گذشت =V+L.گذاشتن
HK1s: کرده باشیم,کرد =V+L.کردن	HP3p: گذاشته اند,گذاشت =V+L.گذاشتن	GP3s: گذاشته بود,گذشت =V+L.گذاشتن
HK3p: کرده باشند,کرد =V+L.کردن	HK2s: گذاشته باشی,گذاشت =V+L.گذاشتن	GP1s: گذاشته بودم,گذشت =V+L.گذاشتن
GP3s: کرده بود,کرد =V+L.کردن	HK1s: گذاشته باشیم,گذاشت =V+L.گذاشتن	GF3s: گذاشته بوده ,گذشت =V+L.گذاشتن
GP2s: کرده بودی,کرد =V+L.کردن	HK3p: گذاشته باشند,گذاشت =V+L.گذاشتن	HP3s: گردونده ,گردوند =V+L.گرداندن

Figure 18. First entries of the compound tenses dictionary

3. Evaluation

Our evaluation was two-fold:

- Evaluate *the formal correction* of the dictionary entries by using the BNFs for both simple and compound forms.
- Evaluate (or rather, estimate) *the lexical coverage* of the dictionaries by applying them to a sample text.

3.1. Evaluation of the formal correction of the dictionary

The inflected simple forms' dictionary contains 1,536 different entries. After applying the simple forms BNF (**Figure 1**, page 8) no errors were found and all entries of the dictionary were matched by the BNF graph.

As for the compound forms retrieved from the corpus, we have 3,953 different entries. After applying the Compound form BNF (**Figure 2**), all entries were also matched, and no errors were found.

3.2. Evaluation of the lexical coverage of the dictionary

A sample text of around 1,000 words was selected from the web. This text is about the world's fastest car that is being built in a workshop on the outskirts of Bristol, England²⁰. (the corresponding translation is also available at the same address). The text was POS-tagged by the system.

²⁰<http://persian.euronews.com/2014/05/01/bloodhound-ssc-the-1000-mph-car-doing-the-school-run>,
<http://www.euronews.com/2014/05/01/bloodhound-ssc-the-1000-mph-car-doing-the-school-run/>
Last access: 02/May/2014

3.2.1. Lexical coverage of simple forms' dictionary

The system found 43 simple verb entries from the simple word dictionary (DLF), corresponding to 107 instances in the text; after manual evaluation, no errors were found in the DLF, and 225 forms were unknown, that is, they are not in the dictionary because they correspond to other POS. These constitute 404 instances in the text (approx. 40%).

Naturally, only the verbs were matched, which means that if a word could belong to other POS, that POS could have been tagged as a verb. For example, the word *آيا* "or" is ambiguous with an inflected form of the verb *آمدن* *amadan* "come" in the simple verb root of spoken form; so naturally, it was matched even though the POS was incorrectly assigned. 21 of 43 simple verb forms were in the same condition. These cases were ignored.

On the other hand, after manual verification, 1 verb from among the unknown words had not been captured by our dictionary. It is the verb *بودن* *bodan* "to be" in third person plural of the spoken form that appears two times in the text. Our dictionary did not recognize it because we forgot to add this inflection of *بود* *bod* "be" to the corresponding inflection graph, which was afterwards corrected.

To sum up²¹, we obtained a Precision of 1, a Recall of 0.98 and F-measure of 0.99.

3.2.2. Lexical coverage of compound forms' dictionary

In the second part of this evaluation, we now deal with compound forms. In this text, we have 31 compound lexical entries and 20 different forms, corresponding to 46 instances in the text (DLC)²².

²¹ Precision (P)= Number of correct verb forms in the text recognized by the dictionary over the total number of verb forms recognized by the dictionary; Recall (R)= number of the correct verb forms in the text recognized by the dictionary over total number of verb forms in the text, and F-measure (F)= $2PR/P+R$

²² For producing a DLC, we apply our compound-tenses graph to the corpus and then we save the resulting concordance as a DLC after some post edition. This includes removing the morphological tags produced for the

In the end, for the existing verbs matched by the systems, no errors were found, which includes the adequate stem, the lemma and the inflection.

As for the evaluation of the compound verbs, we noticed the unknown form *nist* نیست “not to be”, and this form appears only one time and correspond to the contraction of the negation morpheme *ne-* “not” and the third person singular of simple present tense of *هست* *hast* “to be”.

Since this form was missing from the graph for the verb *هست* *hast* “to be”, it was therefore corrected.

To sum up, for the compound forms’ dictionary, we obtained a Precision of 1, Recall of 0.96 and F-measure of 0.98.

Obviously, we are well aware of the size limitation of the sample text, and do not extrapolate these results any further.

In conclusion, though the dictionary is obviously incomplete, since there are many verbs that must still be included, this small experiment provided some evidence that its lexical coverage may be already quite satisfactory. Naturally, in future work, we will endeavor to extend the list of verbs encoded in the DELAS.

simple verb that is part of the compound form: {کردن: HP3p > GS3s:HP3p} > کرد. *کردن*: HP3p (the verb *kard* “do” in the compound present perfect tense *karde and “have been done”*).

4. Conclusion

We have achieved the main objectives of this project. We built a sufficiently large lexicon of Persian verbs, until we attained a frequency threshold corresponding to a cumulative percentage of 90% of the total corpus word forms (3,349,920) corresponding to 1,466 different verb forms. We, therefore, expect to have achieved a reasonable lexical coverage.

We built 9 inflectional graphs to generate automatically all simple forms of Persian verbs; 4 graphs for the written forms and 5 graphs for “spoken” forms, in order to produce a dictionary of simple verb inflected forms, aiming at a highly granular and systematic description of all the relevant grammatical values pertaining to those forms.

Furthermore, we built 22 inflectional graphs to produce automatically all compound forms of Persian verbs, in order to produce a dictionary of compound verb inflected forms.

We produced a Persian verbs dictionary of lemmas (145), where each stem (399) is associated with its corresponding inflectional paradigms; each lemma is associated with 3 different stems (present, past and spoken stem); the distinction between lemmas and stems is overcome by the lemma being treated as a feature of the verb, so one can use it to capture all inflected forms associated with the same lemma with lexical mask in much the same way as one would do for other languages without the distinction between stems and lemmas.

In order to evaluate the correctness of the both dictionaries, we applied the simple form BNF and compound form BNF to the main corpus, and no errors were found in both of them.

In order to evaluate the dictionary lexical coverage, we have applied the dictionaries to a sample of +1000 words text and, for the simple, form we find one verb that appears two times in the text

that was not captured by the dictionary due to forgetting to add this inflection to the DELAS and in the compound form evaluation, after manual verification, we find also one verb that appear one time in the text that was not taken because of the incomplete graph for the verb هست *hast* “to be”.

Naturally, much is still left to be done: First, we want to complete the dictionary by adding many more verbs and, eventually, adding new inflectional paradigms and their corresponding inflection graphs.

Secondly, in order to continue to produce a full lexicon for Persian, attention must be given to another part of speech: this will be the challenge for future work.

References

- Instructional Technology Services group. (2007). Persian online-grammar & resources. Retrieved March 10, 2013, from http://sites.la.utexas.edu/persian_online_resources/verbs/
- Paumier. S. (2003). *De la reconnaissance de formes linguistiques a l'analyse syntaxique*. Univ. Marne-la-Vallée.
- Paumier. S. (2013). *Unitex manual*. Univ. Marne-la-Vallée. Retrieved from <http://www-igm.univ-mlv.fr/~unitex/>
- Persian in Texas group. (n.d.). Persian Simple Verbs. Retrieved December 13, 2013, from <http://persian.nmelrc.org/pvc/simpleverbs.php>
- Persian verbs. (2014). Retrieved June 10, 2013, from http://en.wikipedia.org/wiki/Persian_verbs
- Pilevar, M.T; H. Faili, and A. H. P. (n. d. . (n.d.). Natural Language Processing (NLP) Research lab, University of Tehran. Retrieved May 10, 2013, from <http://ece.ut.ac.ir/NLP/resources.htm>
- Tabatabayi, A. (2010). *ساختمان واژه و مقوله دستوری: تشخیص مقوله دستوری واژه‌ها، بر اساس ملاک‌های صرفی*, *sakhtemane vaje va magholeye dastori: tashkhis magholeye dastoriye vajeha, bar asase melakhaye sarfi, vocabulary structure and grammatical categories: recognition of words grammar* Tahran: Farhang, honar va ertebatat. Retrieved from <http://dlib.ical.ir/site/catalogue/789489>
- فعل ساده, feele sade, Simple verb. (2014). Retrieved November 07, 2013, from http://fa.wikipedia.org/wiki/فعل_ساده

Appendices

A. DELAS (dictionary of Persian verb stems and associated lemmas)

پک, V002+L=پکیدن	گرفت, V002+L=گرفتن	خواب, V001+L=خوابیدن	ماند, V002+L=ماندن
پکوند, V009+L=پکیدن	گرفت, V006+L=گرفتن	خوابید, V002+L=خوابیدن	ماند, V006+L=ماندن
پکید, V002+L=پکیدن	گریخت, V002+L=گریختن	خوابید, V006+L=خوابیدن	مرد, V002+L=مردن
پکید, V006+L=پکیدن	گریخت, V006+L=گریختن	خواست, V002+L=خواستن	مرد, V006+L=مردن
پایید, V002+L=پاییدن	گریز, V001+L=گریختن	خواست, V006+L=خواستن	مون, V008+L=ماندن
پایید, V006+L=پاییدن	گزار, V001+L=گزاردن	خوان, V001+L=خواندن	موند, V002+L=ماندن
پاش, V001+L=پاشیدن	گزارد, V002+L=گزاردن	خواند, V002+L=خواندن	موند, V006+L=ماندن
پاشوند, V009+L=پاشوندن	گزارد, V006+L=گزاردن	خواند, V006+L=خواندن	میر, V001+L=مردن
پاشید, V002+L=پاشیدن	گشت, V002+L=گشتن	خواه, V001+L=خواستن	نگر, V001+L=نگریستن
پاشید, V006+L=پاشیدن	گشت, V006+L=گشتن	خور, V001+L=خوردن	نگریست, V002+L=نگریستن
پخت, V002+L=پختن	گفت, V002+L=گفتن	خور, V008+L=خوردن	نگریست, V006+L=نگریستن
پخت, V006+L=پختن	گفت, V006+L=گفتن	خوران, V001+L=خوراندن	نال, V001+L=نالیدن
پذیر, V001+L=پذیرفتن	گند, V001+L=گندیدن	خوراند, V002+L=خوراندن	نالید, V002+L=نالیدن
پذیرفت, V002+L=پذیرفتن	گندید, V002+L=گندیدن	خورد, V002+L=خوردن	نالید, V006+L=نالیدن
پذیرفت, V006+L=پذیرفتن	گوی, V001+L=گفتن	خورد, V006+L=خوردن	نام, V001+L=نامیدن
پر, V001+L=پریدن	گوز, V001+L=گوزیدن	خورون, V008+L=خوراندن	نامید, V002+L=نامیدن
پراکن, V001+L=پراکندن	گوزید, V002+L=گوزیدن	خوروند, V009+L=خوراندن	نامید, V006+L=نامیدن
پراکند, V002+L=پراکندن	گوزید, V006+L=گوزیدن	داد, V005+L=دادن	نداز, V005+L=انداختن
پراکند, V006+L=پراکندن	گیر, V001+L=گرفتن	داد, V002+L=دادن	نشان, V001+L=نشانندن
پراند, V002+L=پراندن	یا, V001+L=آمدن	داد, V006+L=دادن	نشانند, V002+L=نشانندن
پرداخت, V002+L=پرداختن	یار, V008+L=آوردن	دار, V001+L=داشتن	نشست, V002+L=نشستن
پرداخت, V006+L=پرداختن	آ, V007+L=آمدن	داشت, V002+L=داشتن	نشست, V006+L=نشستن
پرداز, V001+L=پرداختن	آی, V001+L=آمدن	داشت, V006+L=داشتن	نشون, V008+L=نشانندن

نشوند, V009+L=نشانندن	دان, V001+L=دانستن	آشامید, V002+L=آشامیدن	پرس, V001+L=پرسیدن
نشین, V001+L=نشستن	دانست, V002+L=دانستن	آفرید, V002+L=آفریدن	پرسید, V002+L=پرسیدن
نواخت, V002+L=نواختن	دانست, V006+L=دانستن	آفرید, V006+L=آفریدن	پرسید, V006+L=پرسیدن
نواخت, V006+L=نواختن	در, V001+L=دریدن	آفرین, V001+L=آفریدن	پرست, V001+L=پرستیدن
نواز, V001+L=نواختن	درید, V002+L=دریدن	آمد, V002+L=آمدن	پرستید, V002+L=پرستیدن
نوش, V001+L=نوشیدن	درید, V006+L=دریدن	آمد, V006+L=آمدن	پرستید, V006+L=پرستیدن
نوشید, V002+L=نوشیدن	دزد, V001+L=دزدیدن	آموخت, V002+L=آموختن	پرور, V001+L=پروراندن
نوشید, V006+L=نوشیدن	دزدید, V002+L=دزدیدن	آموخت, V006+L=آموختن	پروراند, V002+L=پروراندن
نوشت, V002+L=نوشتن	دزدید, V006+L=دزدیدن	آموز, V001+L=آموختن	پروروند, V009+L=پروراندن
نوشت, V006+L=نوشتن	ده, V001+L=دادن	آمیخت, V002+L=آمیختن	پروند, V009+L=پراندن
نویس, V001+L=نوشتن	دو, V001+L=دویدن	آمیخت, V006+L=آمیختن	پرید, V002+L=پریدن
هست, V004+L=بودن	دوخت, V002+L=دوختن	آمیز, V001+L=آمیختن	پرید, V006+L=پریدن
وز, V001+L=وزیدن	دوخت, V006+L=دوختن	آورد, V001+L=آوردن	پز, V001+L=پختن
وزید, V002+L=وزیدن	دوز, V001+L=دوختن	آورد, V002+L=آوردن	پسند, V001+L=پسندیدن
یاب, V001+L=یافتن	دوش, V001+L=دوشیدن	آورد, V006+L=آوردن	پسندید, V002+L=پسندیدن
یافت, V002+L=یافتن	دوشید, V002+L=دوشیدن	آویخت, V002+L=آویختن	پسندید, V006+L=پسندیدن
یافت, V006+L=یافتن	ساز, V001+L=ساختن	بخش, V002+L=بخشیدن	پندار, V001+L=پنداشتن
فهم, V001+L=فهمیدن	سوخت, V002+L=سوختن	بخشید, V002+L=بخشیدن	پنداشت, V002+L=پنداشتن
فهمید, V002+L=فهمیدن	سوخت, V006+L=سوختن	بخشید, V006+L=بخشیدن	پوک, V008+L=پکیدن
فهمید, V006+L=فهمیدن	سوز, V001+L=سوختن	برد, V001+L=بریدن	پوس, V001+L=پوسیدن
کش, V001+L=کشتن	سوزان, V001+L=سوزاندن	برد, V001+L=بریدن	پوسید, V002+L=پوسیدن
کش, V001+L=کشیدن	سوزاند, V002+L=سوزاندن	برد, V002+L=بردن	پوسید, V006+L=پوسیدن
کشت, V002+L=کشتن	سوزون, V008+L=سوزاندن	برد, V006+L=بردن	پوش, V001+L=پوشیدن
کشت, V006+L=کشتن	سوزوند, V009+L=سوزاندن	برید, V002+L=بریدن	پوشاند, V002+L=پوشاندن
کشید, V002+L=کشیدن	ش, V005+L=شدن	برید, V006+L=بریدن	پوشید, V002+L=پوشیدن
کشید, V006+L=کشیدن	شاش, V001+L=شاشیدن	بست, V002+L=بستن	پوشید, V006+L=پوشیدن

کوب، V001+L=کوبیدن	شاشید، V002+L=شاشیدن	بست، V006+L=بستن	پیچ، V001+L=پیچیدن
کوبید، V002+L=کوبیدن	شتاب، V001+L=شتافتن	بلع، V001+L=بلعیدن	پیچاند، V002+L=پیچیدن
کوبید، V006+L=کوبیدن	شتافت، V002+L=شتافتن	بلعید، V002+L=بلعیدن	پیچوند، V009+L=پیچیدن
کوش، V001+L=کوشیدن	شتافت، V006+L=شتافتن	بلعید، V006+L=بلعیدن	پیچید، V002+L=پیچیدن
کوشید، V002+L=کوشیدن	شد، V002+L=شدن	بند، V001+L=بستن	پیچید، V006+L=پیچیدن
کوشید، V006+L=کوشیدن	شد، V006+L=شدن	بود، V002+L=بودن	پیمای، V001+L=پیمودن
لیس، V001+L=لیسیدن	شست، V002+L=شستن	بود، V006+L=بودن	پیمود، V002+L=پیمودن
لیسید، V002+L=لیسیدن	شست، V006+L=شستن	بوس، V001+L=بوسیدن	پیمود، V006+L=پیمودن
لیسید، V006+L=لیسیدن	شکان، V001+L=شکاندن	بوسید، V002+L=بوسیدن	پیوست، V002+L=پیوستن
مال، V001+L=مالیدن	شکاند، V002+L=شکاندن	بوسید، V006+L=بوسیدن	پیوست، V006+L=پیوستن
مالید، V002+L=مالیدن	شکست، V002+L=شکستن	بوی، V001+L=بوییدن	پیوند، V001+L=پیوستن
مالان، V001+L=مالاندن	شکست، V006+L=شکستن	بویید، V002+L=بوییدن	چپوند، V009+L=چپوندن
مالاند، V002+L=مالاندن	شکن، V001+L=شکستن	بویید، V006+L=بوییدن	چین، V001+L=چیدن
مالون، V008+L=مالاندن	شکون، V008+L=شکاندن	بین، V001+L=دیدن	چاپ، V001+L=چاپیدن
مالوند، V009+L=مالاندن	شکوند، V009+L=شکاندن	تپ، V001+L=تپیدن	چاپید، V002+L=چاپیدن
مان، V001+L=ماندن	شمار، V001+L=شمردن	تپید، V002+L=تپیدن	چاپید، V006+L=چاپیدن
جوشید، V002+L=جوشیدن	شمرد، V002+L=شمردن	تپید، V006+L=تپیدن	چراند، V002+L=چراندن
جوشید، V006+L=جوشیدن	شمرد، V006+L=شمردن	تراش، V001+L=تراشیدن	چرخ، V001+L=چرخیدن
جوی، V001+L=جوییدن	شناخت، V002+L=شناختن	تراشید، V002+L=تراشیدن	چرخاند، V002+L=چرخاندن
جوید، V002+L=جوییدن	شناخت، V006+L=شناختن	تراشید، V006+L=تراشیدن	چرخوند، V009+L=چرخاندن
جوید، V006+L=جوییدن	شناس، V001+L=شناختن	ترس، V001+L=ترسیدن	چرخید، V002+L=چرخیدن
خار، V001+L=خاریدن	شناسان، V001+L=شناساندن	ترسید، V002+L=ترسیدن	چرخید، V006+L=چرخیدن
خاران، V001+L=خارانندن	شناساند، V002+L=شناساندن	ترسید، V006+L=ترسیدن	چروند، V009+L=چراندن
خاراند، V002+L=خارانندن	شناسون، V008+L=شناساندن	تکان، V001+L=تکاندن	چس، V001+L=چسیدن
خارون، V008+L=خارانندن	شناسوند، V009+L=شناساندن	تکاند، V002+L=تکاندن	چسب، V001+L=چسبیدن

خاروند, V009+L=خاراندن	شنو, V001+L=شنیدن	تکاند, V006+L=تکاندن	چسباند, V002+L=چسباندن
خارید, V002+L=خاریدن	شنید, V002+L=شنیدن	توان, V001+L=توانستن	چسیوند, V009+L=چسباندن
خارید, V006+L=خاریدن	شنید, V006+L=شنیدن	توانست, V002+L=توانستن	چسبید, V002+L=چسبیدن
خر, V001+L=خریدن	شو, V001+L=شدن	توانست, V006+L=توانستن	چسبید, V006+L=چسبیدن
خراشید, V002+L=خراشیدن	شوی, V001+L=شستن	تون, V008+L=توانستن	چسید, V002+L=چسیدن
خراشوند, V009+L=خراشانندن	شین, V008+L=نشستن	تونست, V002+L=توانستن	چسید, V006+L=چسیدن
خرید, V002+L=خریدن	صرف, V001+L=صرفیدن	تونست, V006+L=توانستن	چش, V001+L=چشیدن
خرید, V006+L=خریدن	صرفید, V002+L=صرفیدن	جست, V002+L=جستن	چشای, V001+L=چشانندن
خشک, V001+L=خشکیدن	طلب, V001+L=طلبیدن	جست, V006+L=جستن	چشید, V002+L=چشیدن
خشکید, V002+L=خشکیدن	طلبید, V002+L=طلبیدن	جنگ, V001+L=جنگیدن	چشید, V006+L=چشیدن
خشکید, V006+L=خشکیدن	فرست, V001+L=فرستادن	جنگید, V002+L=جنگیدن	چک, V001+L=چکیدن
خفت, V002+L=خفتن	فرستاد, V002+L=فرستادن	جنگید, V006+L=جنگیدن	چکاند, V002+L=چکاندن
خند, V001+L=خندیدن	فرستاد, V006+L=فرستادن	جه, V001+L=جهیدن	چکوند, V009+L=چکاندن
خندید, V002+L=خندیدن	فروخت, V002+L=فروختن	جهید, V002+L=جهیدن	چکید, V002+L=چکیدن
خندید, V006+L=خندیدن	فروخت, V006+L=فروختن	جهید, V006+L=جهیدن	چکید, V006+L=چکیدن
خوا, V007+L=خواستن	فروش, V001+L=فروختن	جوش, V001+L=جوشیدن	چوس, V008+L=چسیدن
گری, V001+L=گریستن	گذار, V001+L=گذاشتن	کرد, V002+L=کردن	چوسید, V009+L=چسیدن
گریست, V002+L=گریستن	گذاشت, V002+L=گذاشتن	کرد, V006+L=کردن	چید, V002+L=چیدن
گرد, V001+L=گذاشتن	گذاشت, V006+L=گذاشتن	کن, V001+L=کردن	چید, V006+L=چیدن
گردان, V001+L=گرداندن	گذر, V001+L=گذشتن	کند, V002+L=کندن	کار, V001+L=کاشتن
گرداند, V002+L=گرداندن	گذشت, V002+L=گذشتن	کند, V006+L=کندن	کاشت, V002+L=کاشتن
گردون, V008+L=گرداندن	گذشت, V006+L=گذشتن	گ, V005+L=گفتن	کاشت, V006+L=کاشتن
گردوند, V009+L=گرداندن			

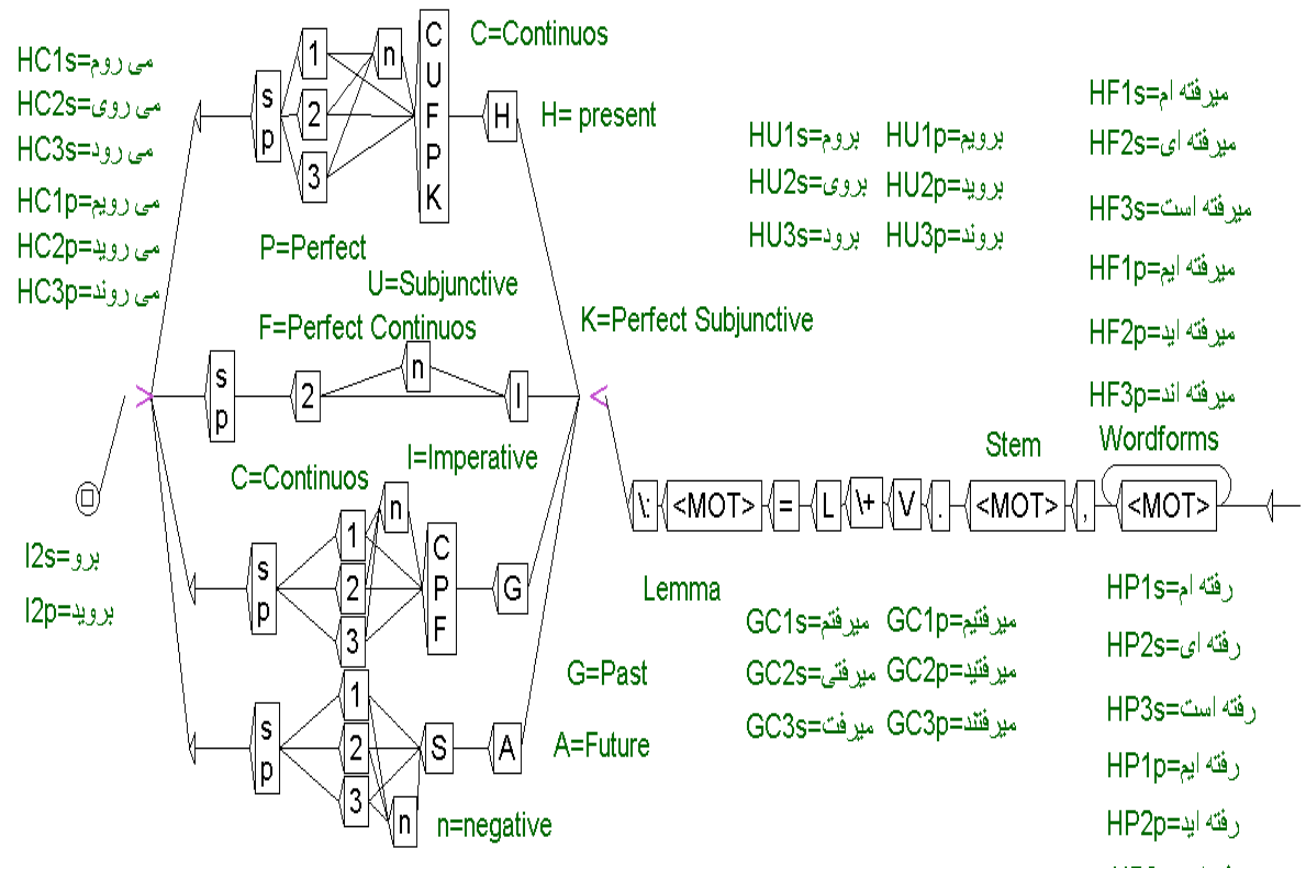
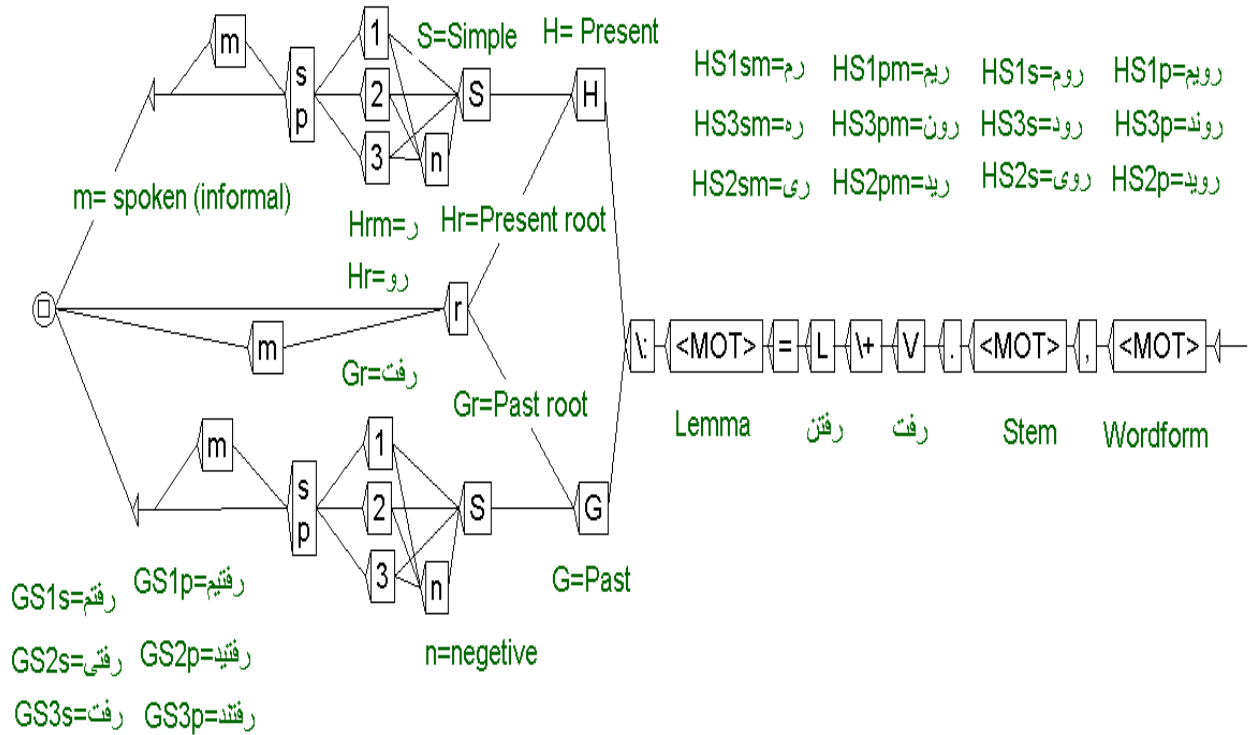
B. Sample of DELAF (dictionary of inflected simple forms)

GS1s: پک، پک=V+L. پکيدن	GS1s: پاييد، پاييد=V+L. پاييدن	GS1s: پاشيد، پاشيد=V+L. پاشيدن
GS2s: پک، پک=V+L. پکيدن	GS2s: پاييد، پاييد=V+L. پاييدن	GS2s: پاشيد، پاشيد=V+L. پاشيدن
GS2s: پک، پک=V+L. پکيدن	GS2s: پاييد، پاييد=V+L. پاييدن	GS2s: پاشيد، پاشيد=V+L. پاشيدن
GS3s: پک، پک=V+L. پکيدن	GS3s: پاييد، پاييد=V+L. پاييدن	GS3s: پاشيد، پاشيد=V+L. پاشيدن
GS1p: پک، پک=V+L. پکيدن	GS1p: پاييد، پاييد=V+L. پاييدن	GS1p: پاشيد، پاشيد=V+L. پاشيدن
GS2p: پک، پک=V+L. پکيدن	GS2p: پاييد، پاييد=V+L. پاييدن	GS2p: پاشيد، پاشيد=V+L. پاشيدن
GS3p: پک، پک=V+L. پکيدن	GS3p: پاييد، پاييد=V+L. پاييدن	GS3p: پاشيد، پاشيد=V+L. پاشيدن
Gr: پک، پک=V+L. پکيدن	Gr: پاييد، پاييد=V+L. پاييدن	Gr: پاشيد، پاشيد=V+L. پاشيدن
GS1sm: پکوند، پکوند=V+L. پکيدن	GS2pm: پاييد، پاييد=V+L. پاييدن	GS2pm: پاشيد، پاشيد=V+L. پاشيدن
GS2sm: پکوند، پکوند=V+L. پکيدن	GS3pm: پاييد، پاييد=V+L. پاييدن	GS3pm: پاشيد، پاشيد=V+L. پاشيدن
GS2sm: پکوند، پکوند=V+L. پکيدن	HS1s: پاش، پاش=V+L. پاشيدن	GS1s: پخت، پخت=V+L. پختن
GS3sm: پکوند، پکوند=V+L. پکيدن	HS2s: پاش، پاش=V+L. پاشيدن	GS2s: پخت، پخت=V+L. پختن
GS1pm: پکوند، پکوند=V+L. پکيدن	HS2s: پاش، پاش=V+L. پاشيدن	GS2s: پخت، پخت=V+L. پختن
GS2pm: پکوند، پکوند=V+L. پکيدن	HS3s: پاش، پاش=V+L. پاشيدن	GS3s: پخت، پخت=V+L. پختن
GS2pm: پکوند، پکوند=V+L. پکيدن	HS1p: پاش، پاش=V+L. پاشيدن	GS1p: پخت، پخت=V+L. پختن
GS3pm: پکوند، پکوند=V+L. پکيدن	HS2p: پاش، پاش=V+L. پاشيدن	GS2p: پخت، پخت=V+L. پختن
GS3pm: پکوند، پکوند=V+L. پکيدن	HS3p: پاش، پاش=V+L. پاشيدن	GS3p: پخت، پخت=V+L. پختن
Grm: پکوند، پکوند=V+L. پکيدن	Hr: پاش، پاش=V+L. پاشيدن	Gr: پخت، پخت=V+L. پختن
GS1s: پکيد، پکيد=V+L. پکيدن	GS1sm: پاشوند، پاشوند=V+L. پاشوندن	GS2pm: پخت، پخت=V+L. پختن
GS2s: پکيد، پکيد=V+L. پکيدن	GS2sm: پاشوند، پاشوند=V+L. پاشوندن	GS3pm: پخت، پخت=V+L. پختن
GS2s: پکيد، پکيد=V+L. پکيدن	GS2sm: پاشوند، پاشوند=V+L. پاشوندن	HS1s: پذير، پذير=V+L. پذيرفتن
GS3s: پکيد، پکيد=V+L. پکيدن	GS3sm: پاشوند، پاشوند=V+L. پاشوندن	HS2s: پذير، پذير=V+L. پذيرفتن
GS1p: پکيد، پکيد=V+L. پکيدن	GS1pm: پاشوند، پاشوند=V+L. پاشوندن	HS2s: پذير، پذير=V+L. پذيرفتن
GS2p: پکيد، پکيد=V+L. پکيدن	GS2pm: پاشوند، پاشوند=V+L. پاشوندن	HS3s: پذير، پذير=V+L. پذيرفتن
GS3p: پکيد، پکيد=V+L. پکيدن	GS2pm: پاشوند، پاشوند=V+L. پاشوندن	HS1p: پذير، پذير=V+L. پذيرفتن
Gr: پکيد، پکيد=V+L. پکيدن	GS3pm: پاشوند، پاشوند=V+L. پاشوندن	HS2p: پذير، پذير=V+L. پذيرفتن
GS2pm: پکيد، پکيد=V+L. پکيدن	GS3pm: پاشوند، پاشوند=V+L. پاشوندن	HS3p: پذير، پذير=V+L. پذيرفتن
GS3pm: پکيد، پکيد=V+L. پکيدن	Grm: پاشوند، پاشوند=V+L. پاشوندن	Hr: پذير، پذير=V+L. پذيرفتن

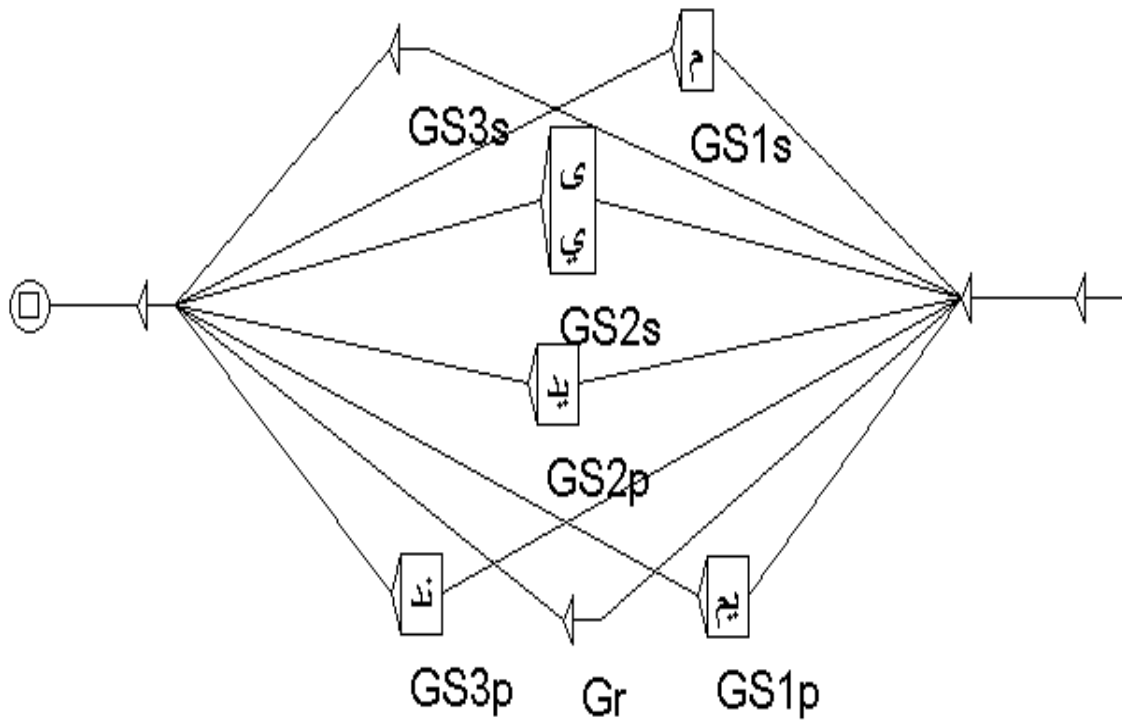
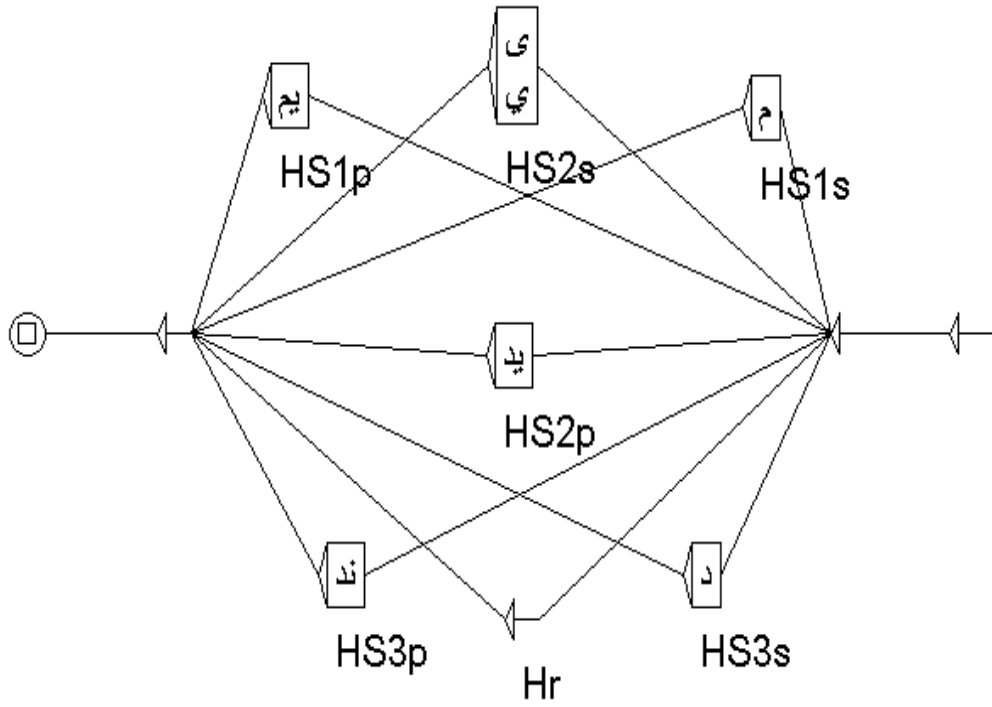
C. Sample of Compound forms (found in a text)

HP3s: پخته=V+L. پخت.	HP3s: پوشیده است، پوشید.	HP3p: پیوسته اند، پیوست.
HP3p: پخته اند، پخت.	HP3p: پوشیده اند، پوشید.	HP3s: چپونده، چپوند.
GP3s: پخته بود، پخت.	HP3s: پوشانده، پوشاند.	GP1s: چپونده بودم، چپوند.
HP3s: پذیرفته، پذیرفت.	HP3s: پوشانده است، پوشاند.	HP3s: چاپیده، چاپید.
HP1s: پذیرفته ام، پذیرفت.	HP3p: پوشانده اند، پوشاند.	HP3s: چرخانده، چرخاند.
HP3p: پذیرفته اند، پذیرفت.	GP3s: پوشانده بود، پوشاند.	HP3s: چرخانده است، چرخاند.
GP3p: پذیرفته بودند، پذیرفت.	GP1s: پوشانده بودم، پوشاند.	HP3s: چرخیده، چرخید.
HP3s: پراکنده، پراکند.	GF3s: پوشانده بوده، پوشاند.	HP3p: چرخیده اند، چرخید.
HP3s: پرداخته، پرداخت.	HP3s: پوشیده، پوشید.	HP3s: چسبانده، چسباند.
HP3s: پرسیده، پرسید.	HP1s: پوشیده ام، پوشید.	GP3s: چسبانده بود، چسباند.
HP1p: پرسیده ایم، پرسید.	HP3p: پوشیده اند، پوشید.	HP3s: چسبونده، چسبوند.
HP3s: پرستیده، پرستید.	HK1s: پوشیده باشم، پوشید.	GP3s: چسبونده بود، چسبوند.
HP3s: پرورانده، پروراند.	GP3s: پوشیده بود، پوشید.	HP3s: چسبیده، چسبید.
GP3p: پرورانده بودند، پروراند.	GP1s: پوشیده بودم، پوشید.	HP3p: چسبیده اند، چسبید.
HP3s: پرونده، پروند.	GP3p: پوشیده بودند، پوشید.	GP3s: چسبیده بود، چسبید.
HP3s: پرونده است، پروند.	GF3s: پوشیده بوده، پوشید.	GP1s: چسبیده بودم، چسبید.
HP1s: پرونده ام، پروند.	HP3s: پیچیده، پیچوند.	HP3s: چشیده، چشید.
GP3s: پرونده بود، پروند.	HK1s: پیچونده باشم، پیچوند.	HP3s: چکیده، چکید.
GF3s: پرونده بوده، پروند.	HP3s: پیچیده، پیچید.	HP3s: چیده، چید.
HP3s: پریده، پرید.	HP3s: پیچیده است، پیچید.	HK3p: چیده باشند، چید.
HP3s: پریده است، پرید.	GP3s: پیچیده بود، پیچید.	HP3s: کاشته، کاشت.
GP3s: پریده بود، پرید.	HP3s: پیموده، پیمود.	GP3s: کاشته بود، کاشت.
HP3s: پسندیده، پسندید.	HP3s: پیوسته، پیوست.	HP3s: کرده، کرد.
HP3s: پنداشته، پنداشت.	HP3s: پیوسته است، پیوست.	HP2s: کرده ای، کرد.
HP3s: پوشیده، پوشید.	HP1s: پیوسته ام، پیوست.	HP2p: کرده اید، کرد.

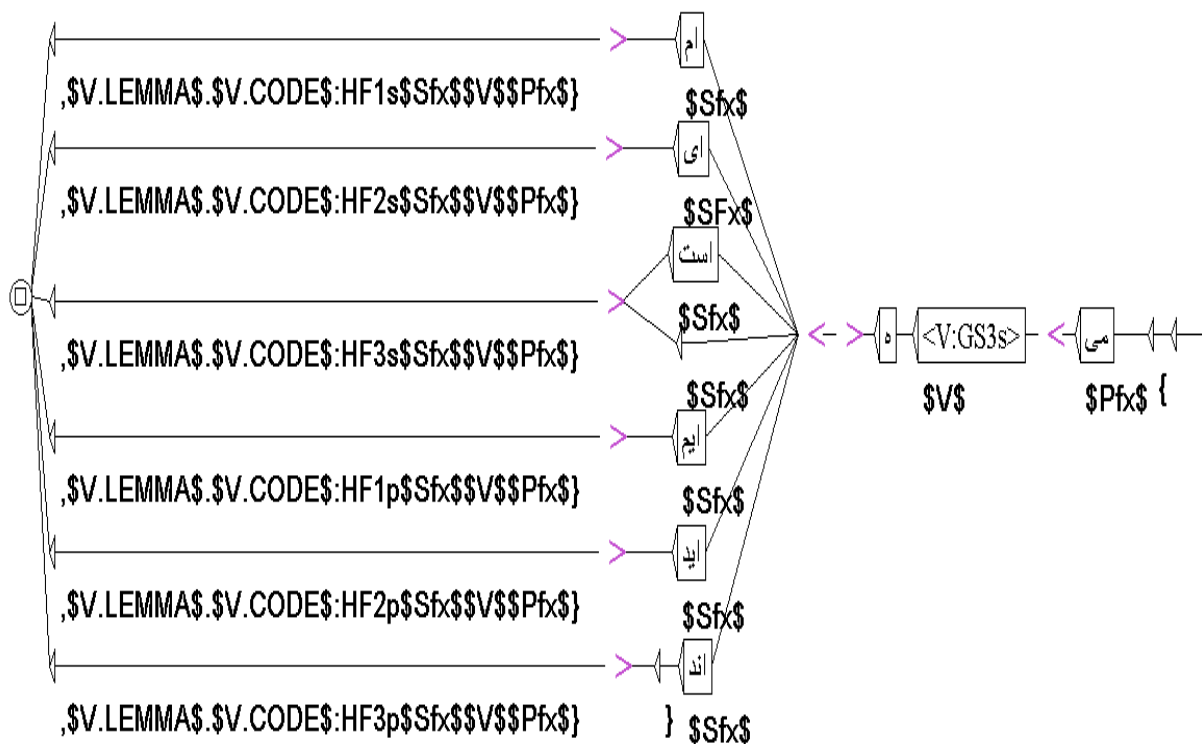
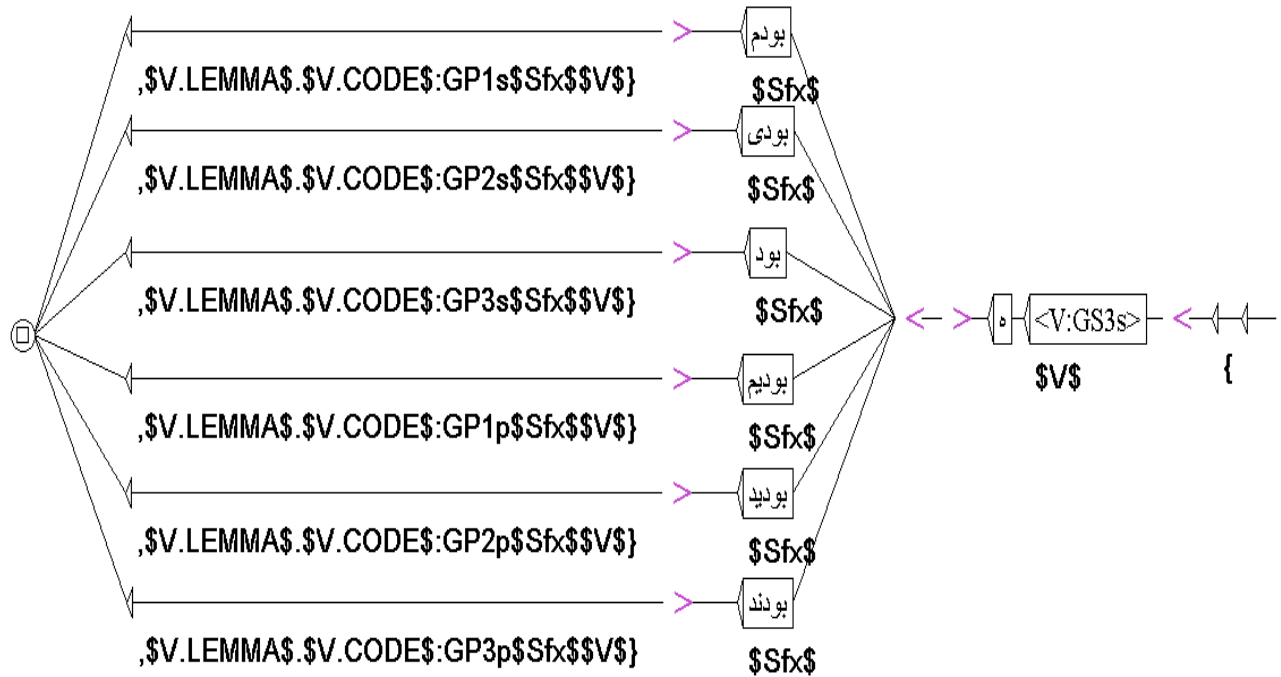
D. Simple BNF and Compound BNF



E. Sample of simple verbs FST



F. Sample of Compound verbs FST



G. Evaluation sample text

In Persian²³:

رازهای غیر پنهان اتومبیلی با سرعت یکهزار مایل در ساعت سریع ترین اتومبیل جهان قطعه قطعه در این کارگاه در بریتانیا ساخته می شود. هدف دست یابی به سرعت ۱۰۰۰ مایل یا ۱۶۰۹ کیلومتر در ساعت است.

این خودروی فراصوتی "بلدهوند اس.اس.سی" نام دارد و پشت فرمان نمونه آزمایشی آن اندی گرین، خلبان نظامی و صاحب رکورد جهانی سرعت نشسته است.

او برای نخستین بار دریچه خودرو را باز می کند تا صندلی راننده را به ما نشان دهد. اندی گرین می گوید: "این اتاقک خودروی فراصوتی بلدهوند اس.اس.سی، یعنی دفتر کار ۱۰۰۰ مایل در ساعت من است. دارای فرمانی با چاپ سه بعدی از تیتانیوم است که با دست های من جور است. همچنین دارای سه صفحه نمایش از کریستال مایع است. دو ابزار محاسبه آنرا رولکس ساخته، یعنی سرعت سنج و ساعت همراه با کرومومتر. دقت این ابزار لازمه کاری است که با خودرو انجام می دهیم"

این مدل، شبیه اتومبیلی است که تا یکسال دیگر تکمیل می شود. و این چیزی است که اکنون در حال ساخته شدن است.

سه موتور روی خودرو نصب است، یک موتور جت، یک موشک و یک موتور خودرو. فولادی بودن بدنه خود رو به خاطر محکم و قوی بودن آنست.

این خودرو یک رویای مهندسی است و رویای تیم سازنده این است که الهام بخش مهندسان جوان باشد. به همین دلیل هم هیچ چیز سری در اینجا وجود ندارد.

گزارشگر یورونیوز می گوید: "یکی از ویژگیهای این پروژه این است که تمام اطلاعات، و تمام طرحهای مربوط به ساخت خودرو و طرز رانندگی آن، برای تمام دنیا آنلاین و آزاد است."

در کنار تلاشی که برای ارائه باز منابع اطلاعات این خودرو می شود، تیم سازنده بلاد هوند هر سال از صدها مدرسه دیدن می کنند. این شیوه تشویق به یاد گیری، از ریچارد نوبل صاحب پیشین رکورد سرعت زمینی است.

ریچارد نوبل می گوید: "بریتانیایی ها در سال به صد هزار مهندس کارشناس نیاز دارند و این برای آماده کردن سی هزار نفر آن ها کافی ست. به این ترتیب وضعیت خیلی جدی است. مساله این است که باید راهی به دبستان و شگفت زده کردن کودکان یافت تا در سالهای بعد جوانان علاقمند شوند."

خودرویی که موتور آن ۱۳۵ هزار اسب قدرت دارد و می تواند سریع تر از جت جنگی حرکت کند، واقعا باید خیلی قابل توجه باشد. و هیچ کس بیشتر از مردی که پشت فرمان آن نشسته توجه برانگیز نیست.

اندی گرین، راننده نمونه آزمایشی این خودرو: "در حال حاضر من یاد می گیرم با تیمی که رکورد بیشترین سرعت زمینی را در تاریخ دارند، کار کنم، مهندسانی در سطح جهانی که مسائل سختی را حل کرده اند. در سال های آینده ما این ماشین را توسعه داده و تکمیل میکنیم و بعد به همراه تیم آن را می رانیم و تمرین می کنیم. همزمان

²³ <http://persian.euronews.com/2014/05/01/bloodhound-ssc-the-1000-mph-car-doing-the-school-run>
Date of access: 02/May/2014

ما جراجویی های خودمان را رسانه ای می کنیم تا از طریق مخاطبان جهانی الهام بخش نسل آینده باشیم.

در طول یکسال آینده بلد هوند اس اس سی آماده می شود تا به آفریقای جنوبی برود و هدف خود را که طی کردن یک هزار مایل در ساعت است، آغاز کند.

In English²⁴:

<Inside Bloodhound SSC: the 1000 mph car>

Piece by piece, the world's fastest car is being built in a workshop on the outskirts of Bristol, England.

Known as Bloodhound SSC, it has the goal of hitting the heady target of 1,000 miles per hour, and inspiring millions of school children in the process.

Behind the wheel of the arrow-shaped machine is current world land speed record holder and military pilot, Andy Green.

For the first time, he opened the hatch to show us the driver's seat.

"This is the cockpit for the Bloodhound Super Sonic Car. This is my 1,000 mile per hour office," Green smiles.

"It has a 3D steering wheel printed in titanium to match my hands.

"We have three liquid crystal displays, and the two backup instruments that Rolex has produced, which is the backup speedometer, and the clock and stopwatch."

"The precision of these instruments is essential to what we're doing," he stresses.

At the other side of the workshop lies a full-scale mock-up of the car, which gives a sense of its scale. At 13.5 metres long and three metres high, it looks like a cross between a jet fighter and a 9-year-old's sketch of their dream dragster.

Bloodhound SSC chasis

Inside there will be more than a few aerospace ingredients in the Bloodhound recipe, including at its heart a Rolls Royce EJ220 jet engine from a Eurofighter Typhoon. This is paired with a Falcon hybrid rocket from Norwegian company Nammo, a supplier to the European space industry.

For the moment the real car is still being built. Chief Engineer Mark Chapman showed us around: "This is the front of the car, this is where the wheels attach to, the front suspension," he explains. "This is all carbon fibre, this is Andy's office," he says, gesturing to the cockpit. The lower assembly behind Green is made from steel. "The reason we've got a steel structure is it's very, very stiff and very very strong," Chapman says.

The rear part of the car where the rear wheels will be mounted is aluminum. "It's very complicated machining, but a beautiful part," he adds proudly.

The car is an engineer's dream, and the dream of the team is to inspire young engineers.

²⁴ <http://www.euronews.com/2014/05/01/bloodhound-ssc-the-1000-mph-car-doing-the-school-run/>
Date of access: 02/May/2014

For that reason, nothing here is secret. All of the data on how the car is built and how it will be tested and driven as it attempts to smash the current world land speed record will be shared with the world online, rights-free.

Alongside that push to offer open source information, the Bloodhound team visit hundreds of schools in Britain every year.

The vision for education came from former land speed record holder and Bloodhound SSC Project Director, Richard Noble.

“Britain needs 100,000 graduate engineers per annum, and the system is only delivering 30,000, so it’s a really serious situation. And the problem actually boils right the way down to primary school, the need to get the kids really excited at a really interesting and young age,” he says.

The team have put enormous effort into creating eye-catching animations to bring the project to life. These include Bloodhound SSC racing and beating a fighter jet, and driving faster than a bullet from a pistol.

In Andy Green they have one of the world’s most convincing ambassadors for the thrill of driving faster than anyone has ever driven before in a car with the equivalent of 135,000 horsepower.

His eyes filled with excitement, he tells us: “right now I get to work with the best land speed record team in history, world class engineers solving problems nobody has solved before. Next year and the year after we actually get to test and develop this car, and I get to drive it and be the delivery part of that test team. And at the same time we get to transmit that message and share the adventure live to a global audience, we get a chance to inspire that next generation.”

The car should be ready for its first test run in spring next year – a humble 200 miles per hour (322 kph) on a Cornish air field. Then in summer next year, the team will decamp to Hakskeen Pan in South Africa, a dry lake that has been specially cleared of small stones to provide a perfect surface for the record attempt.

The aim is to break the current world land speed record of 763 miles per hour, or 1,227 kilometers an hour, in 2015.

Then, in 2016, the car will be readied for its ultimate challenge of breaking through the mythical barrier of 1,000 mph.