# An Electronic Dictionary of Persian Verbs

## Bahareh Kakanaeeni

3rd UNITEX/GramLab Workshop

October 9-10, 2014

# Outline

- Motivation and problem
- Related work
- Objectives
- Methods
- Evaluation
- Conclusion
- future work
- References

# 1.Motivation and Problem



There are around 110 million Persian speakers worldwide. Persian is spoken in Iran, Afghanistan, Tajikistan and other countries that historically came under Persian influence.

Source: http://en.wikipedia.org/wiki/Persian_language

- In spite of its cultural and demographic importance, studies on Persian in view of Natural Language Processing (NLP) are scarce. One of the main problems is the fact that lexical resources for NLP are almost non-existent in Persian.

- Because of this, I intend to start to build a new module for the Unitex linguistic development platform (Paumier 2003, 2013) by building linguistic resources for Persian contribute to lay the basis for further studies of this language in a computational linguistics perspective. I start by describing the inflection of Persian verbs, which is the topic of this presentation.

# 2.Related work

1. Finite-State Morphological Analysis of Persian by Karine Megerdoomian from University of California, San Diego.

2. A Morphological Lexicon for the Persian Language by Benoˆıt Sagot, G´eraldine Walther from Universit´e Paris.

3. Persian Morphology by John R. Perry from University of Chicago.

# 3.Objectives

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

1. Build a sufficiently large lexicon of Persian verbs 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 simple forms of Persian verbs, in order to produce a dictionary of inflected forms and build morphological FST to recognize compound verb forms in the text.

4. Build and apply BNFs to validate the formal correction of the simple words' dictionary and the compound forms matched in texts.

5. Apply the Persian dictionary to a sample of text in order to assess the lexical coverage and grammatical adequacy of the dictionary.

# 4.Methods

## 4.1. Synopsis of Persian verb morphology

| Verb forms | Prefix | Verb stem | Suffix |
|---|---|---|---|
| Simple present (HS) | ___ | rav (ro) "go" | i |
| Simple past (GS) | ___ | raft "gone" | i |
| Past Continuous (GC) | mi | raft "gone" | i |
| Present Continuous (HC) | mi | rav (ro) "go" | i |
| Present Subjunctive (HU) | be | rav (ro) "go" | i |
| Imperative (I) | bo | rav (ro) "go" | ___ |
| Negative (n) | na | rav (ro) "go" | ___ |
| Present Perfect (HP) | ___ | raft "gone" | eh i |
| Past Perfect (GP) | ___ | raft "gone" | eh bodeh i |

## 4.2. A frequency-based lexicon of Persian verbs

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

Source for token frequency values: TEP corpus (4,485,147 words)
[Ref:http://ece.ut.ac.ir/nlp/resources.htm]

## 4.3. A Persian verbs dictionary of inflected forms

## 4.3.1. BNF

BNFs in a form of FSA were built to provide formal guidelines for the description of the verb dictionary entries (simple and compound) and, I use these BNFs in order to validate automatically the inflected words' dictionary entries.

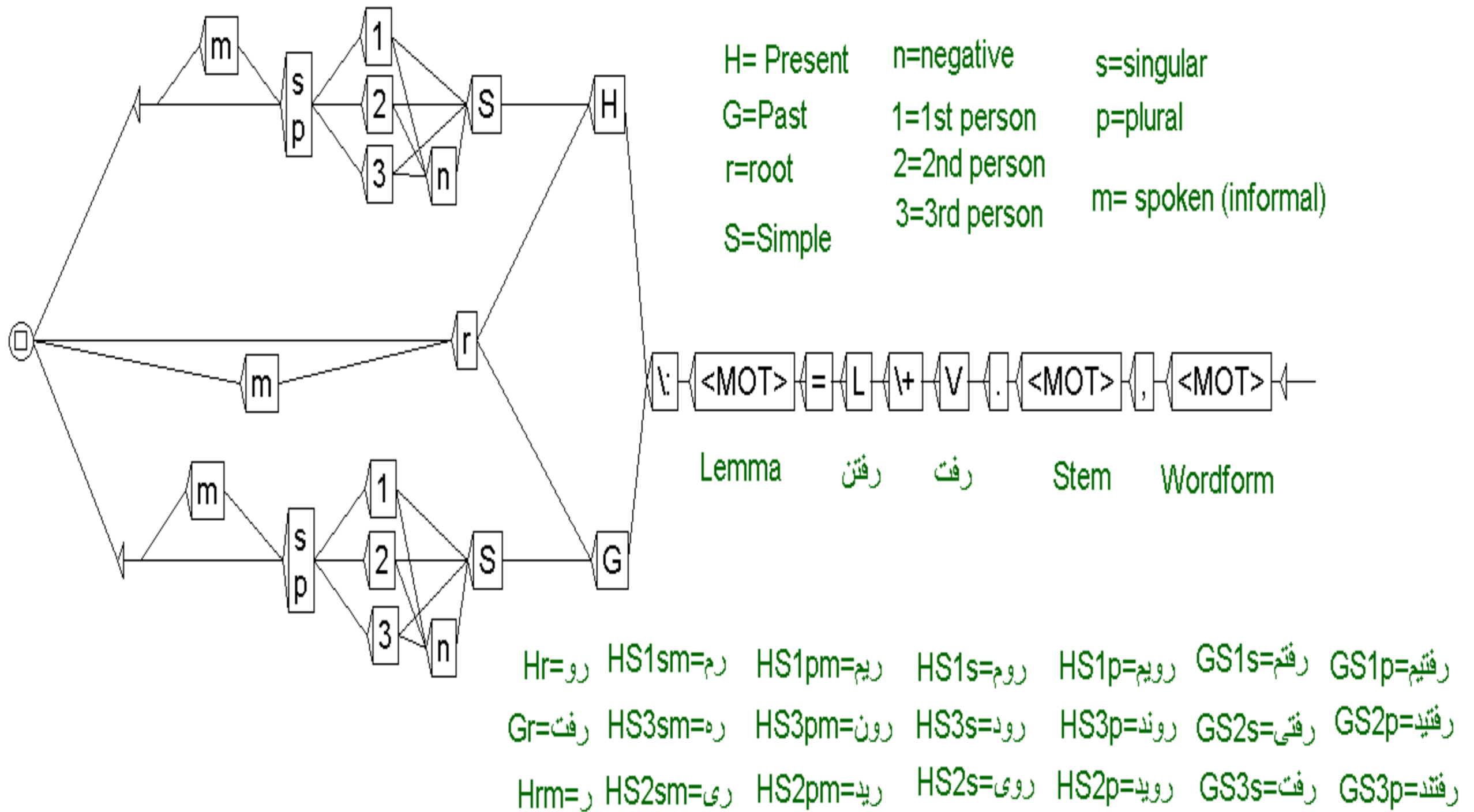Source: http://en.wikipedia.org/wiki/Backus%E2%80%93Naur_Form
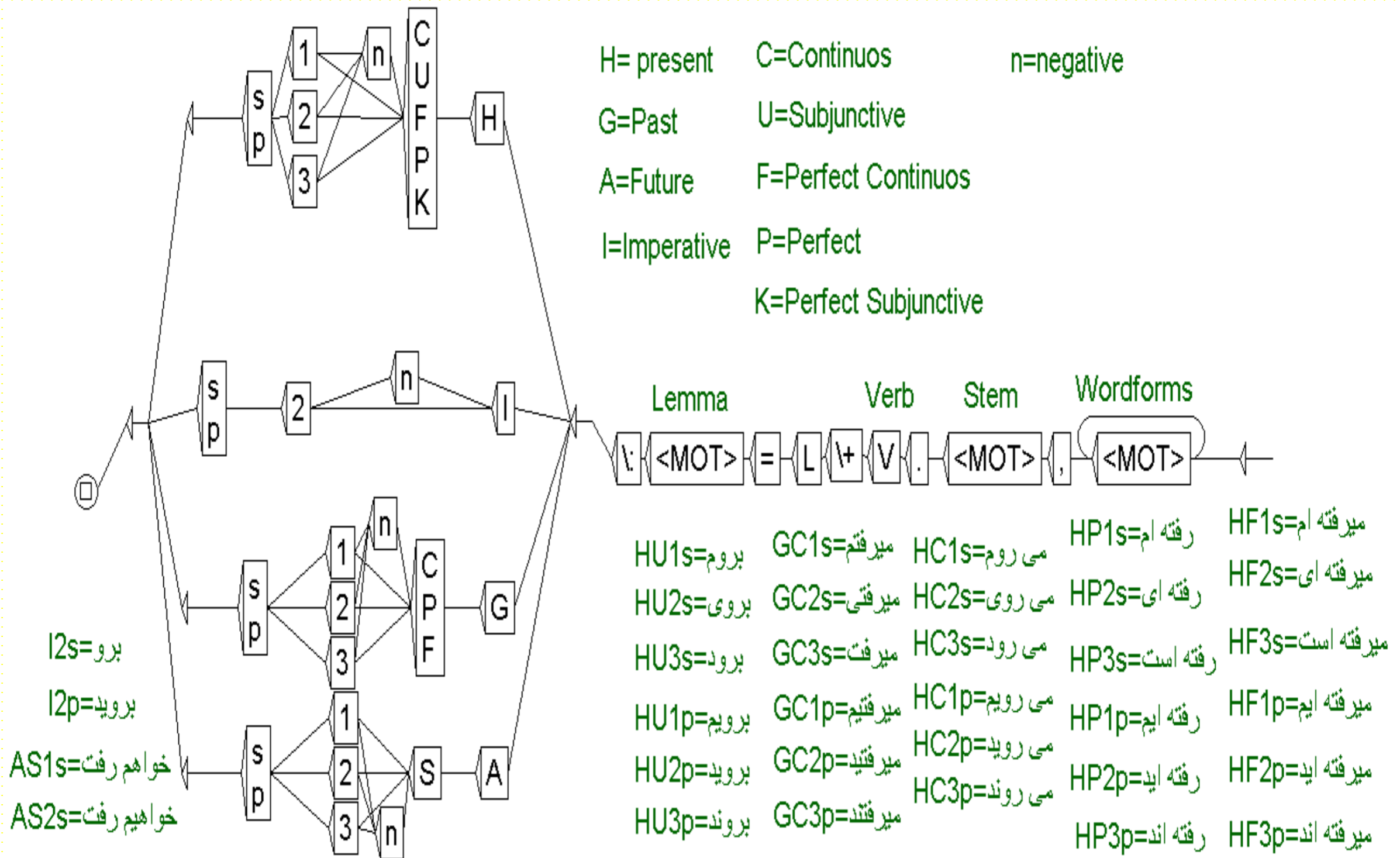
Figure 1. BNF1 for Simple Tenses

Figure 2. BNF2 for Compound Tenses
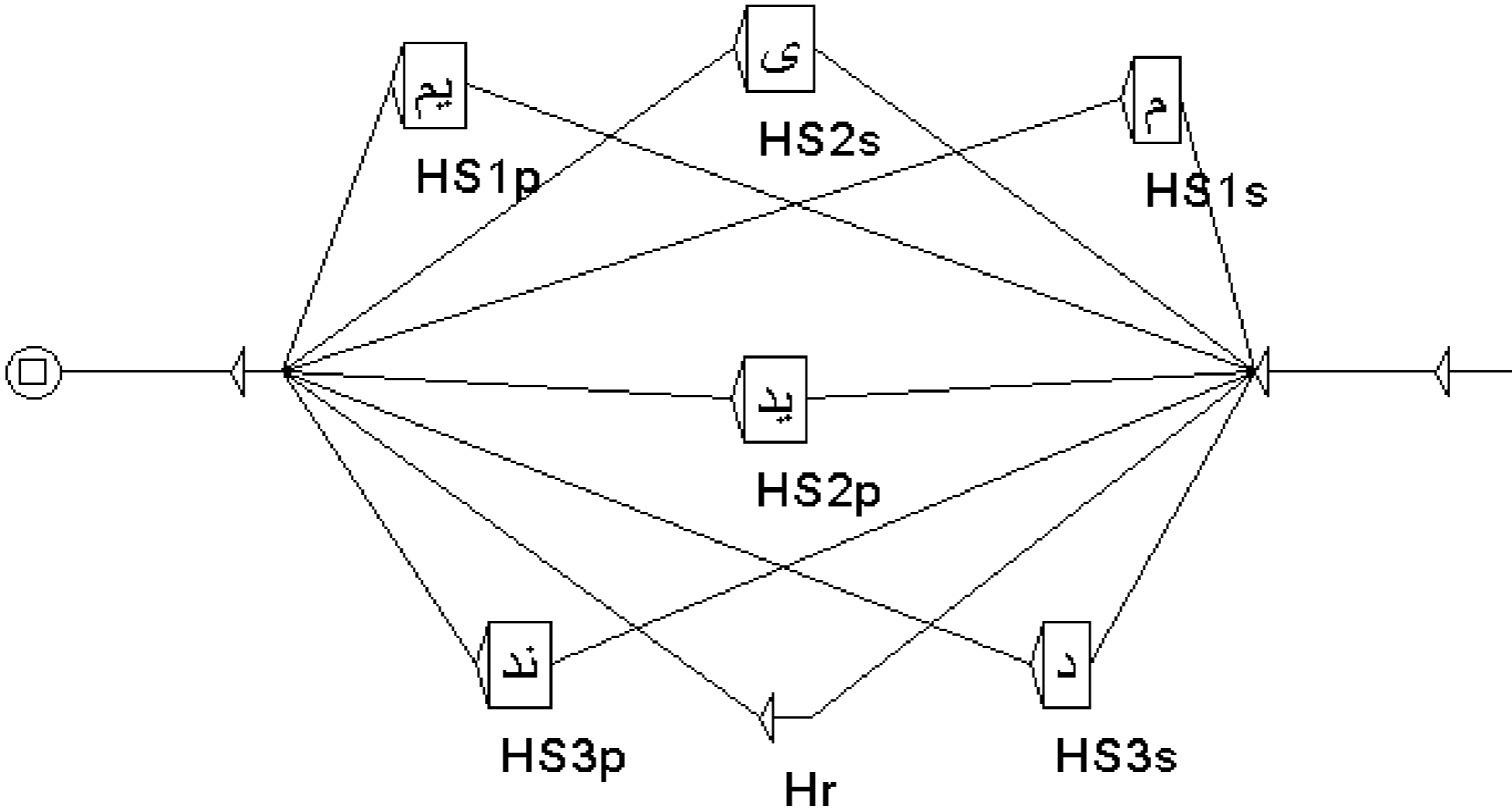
12

# 4.3.2. Inflection FSTs for written simple words



Figure 3. V001, Written Simple Present graph
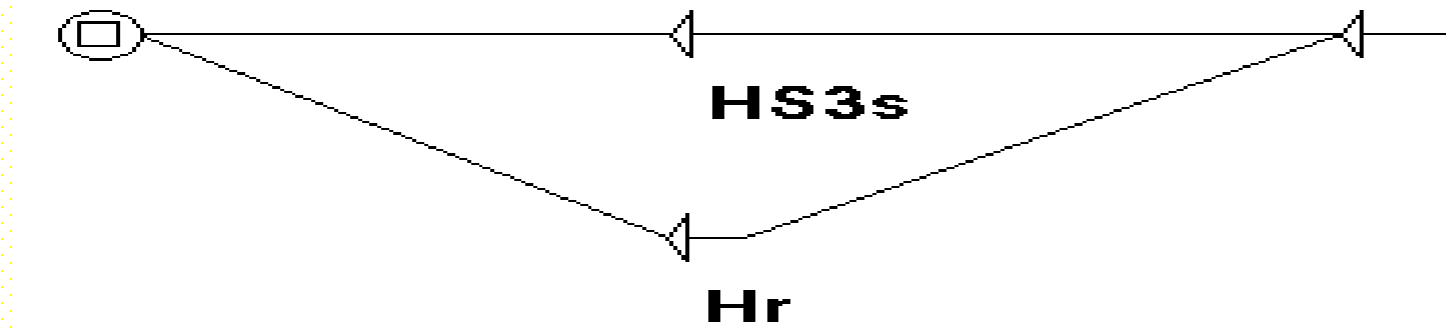
Figure 4. V002, Written Simple Past graph
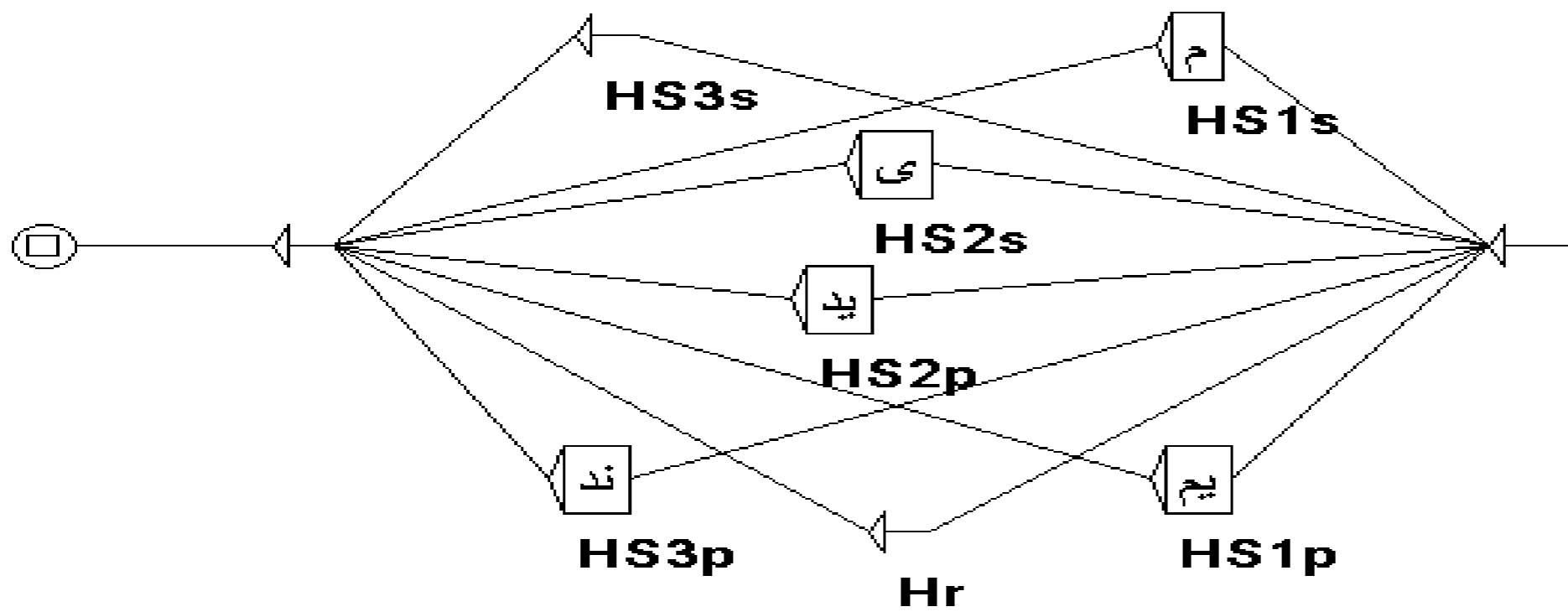
14

Figure 5. V003, است *ast* (be) graph
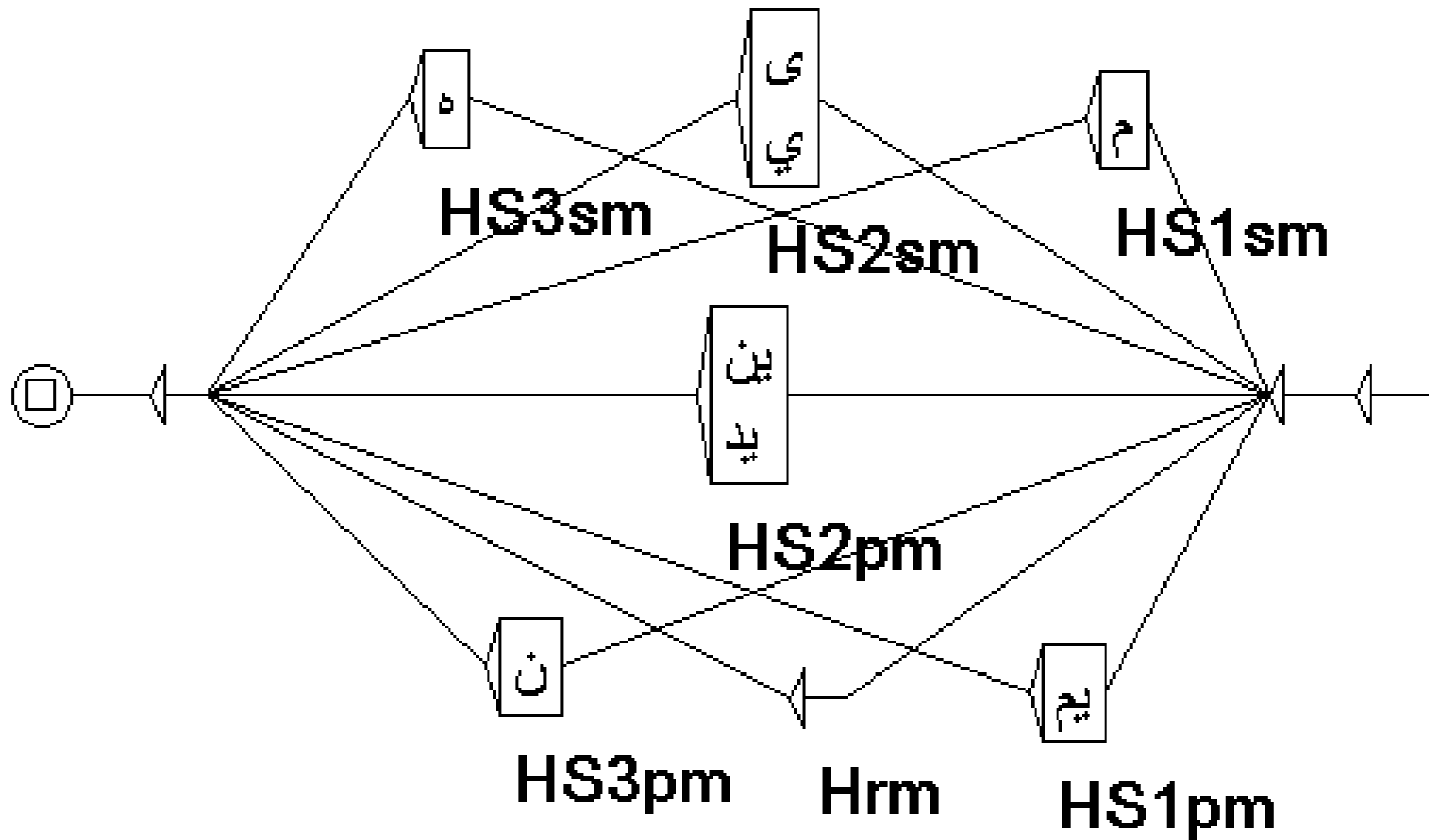


Figure 6. V004, هست *hast* (be) graph

15

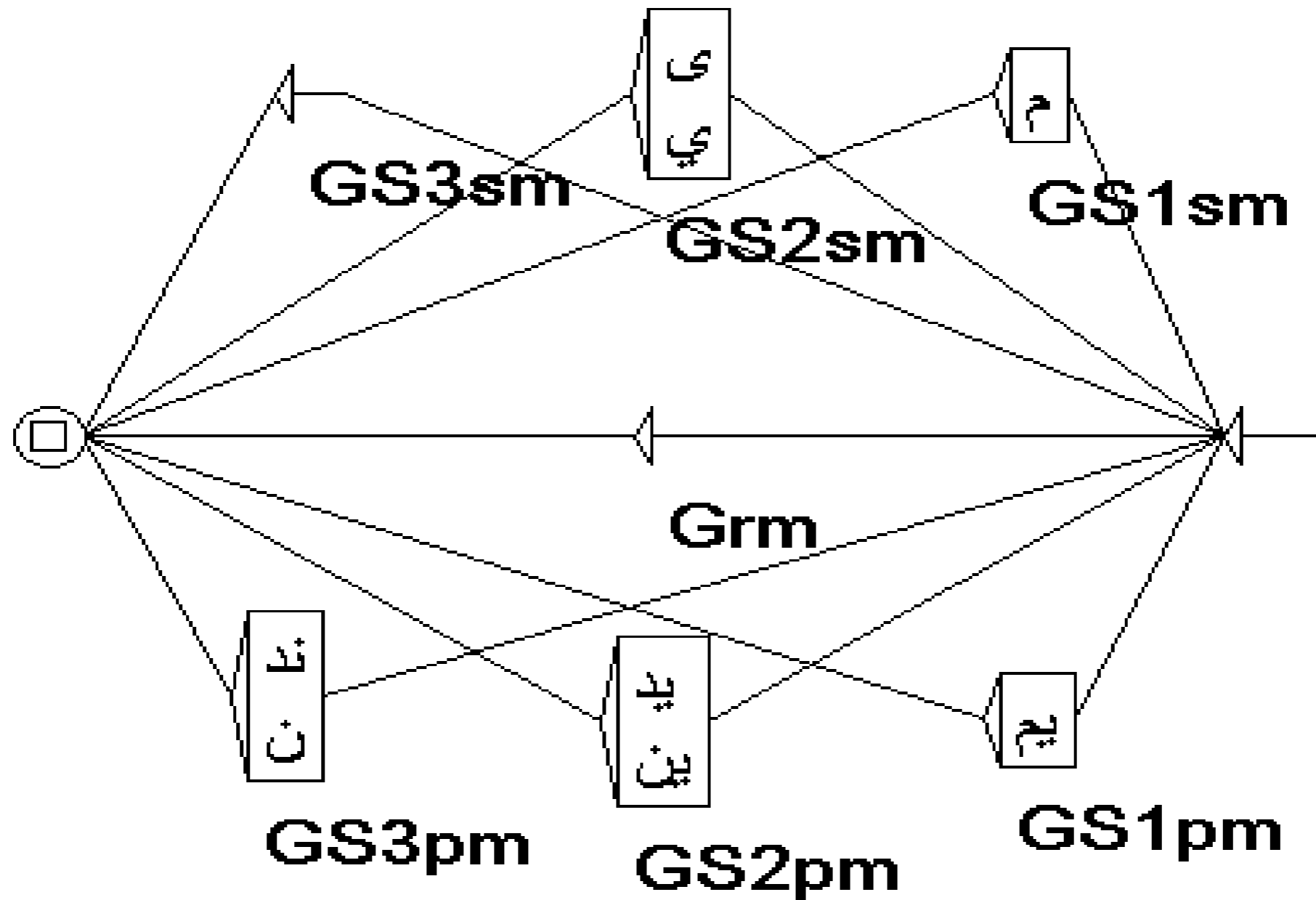Figure 7. V005, Spoken Simple Present graph

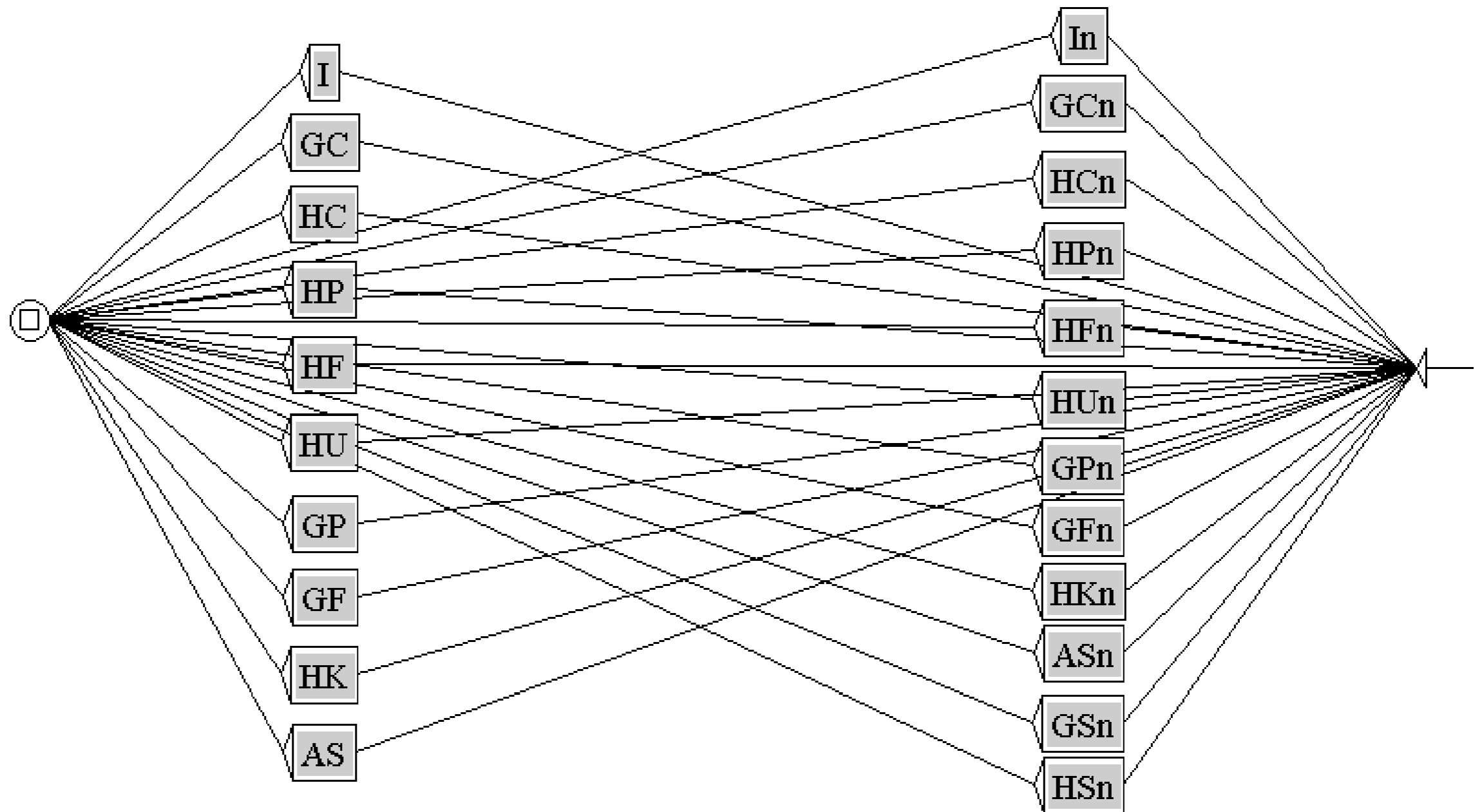Figure 8. V009, Spoken Simple Past graph

# 4.4. Compound inflection FSTs



Figure 9. Compound-tenses general graph

This compound tense is formed by prefix بِ "*be*" and the simple present forms of the verb.
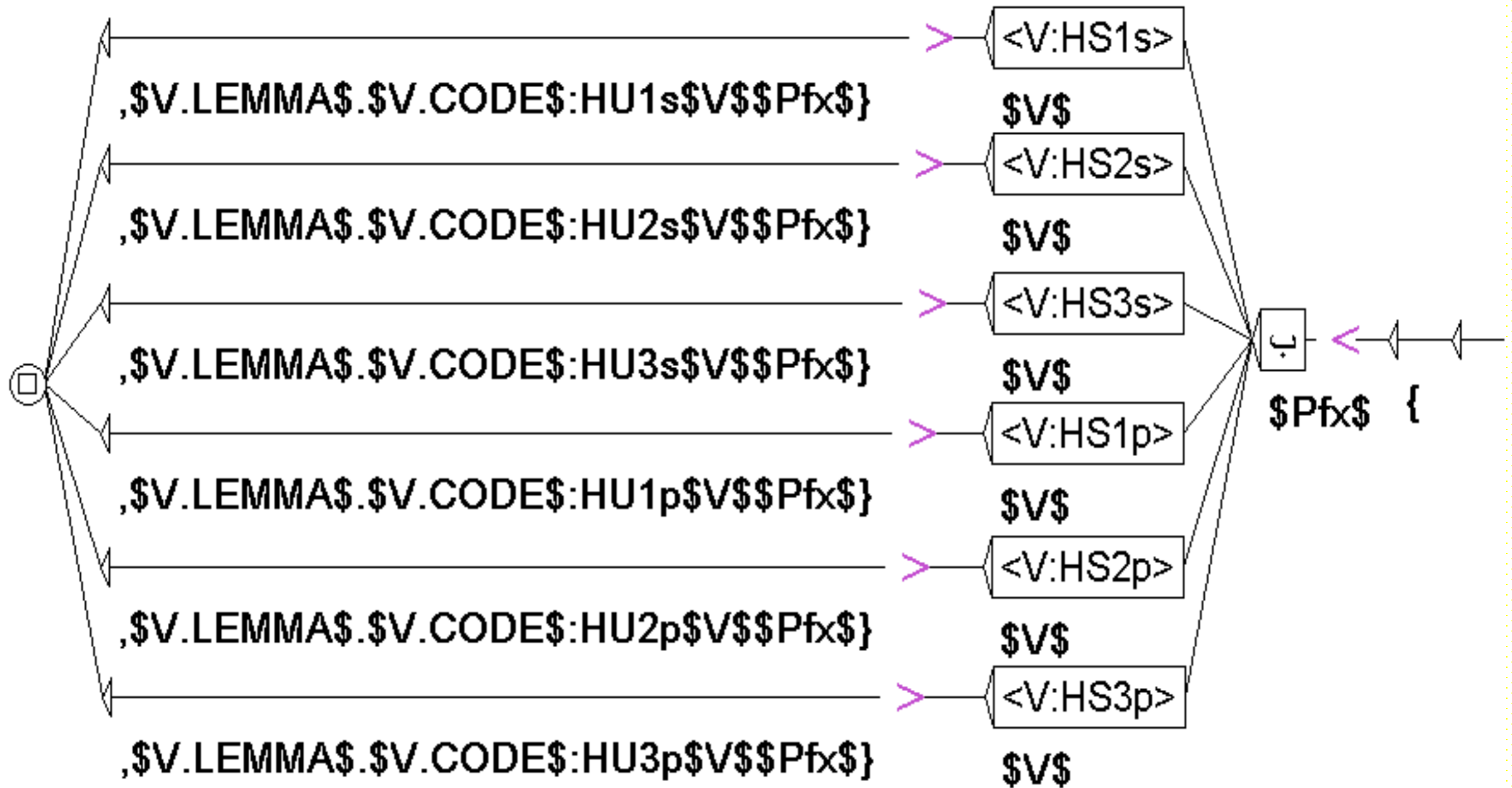


Figure 10. Morphological graph for Present Subjunctive (HU)

In this graph, prefix می "mi" can be a separate token or be joined with the base form of the verb.
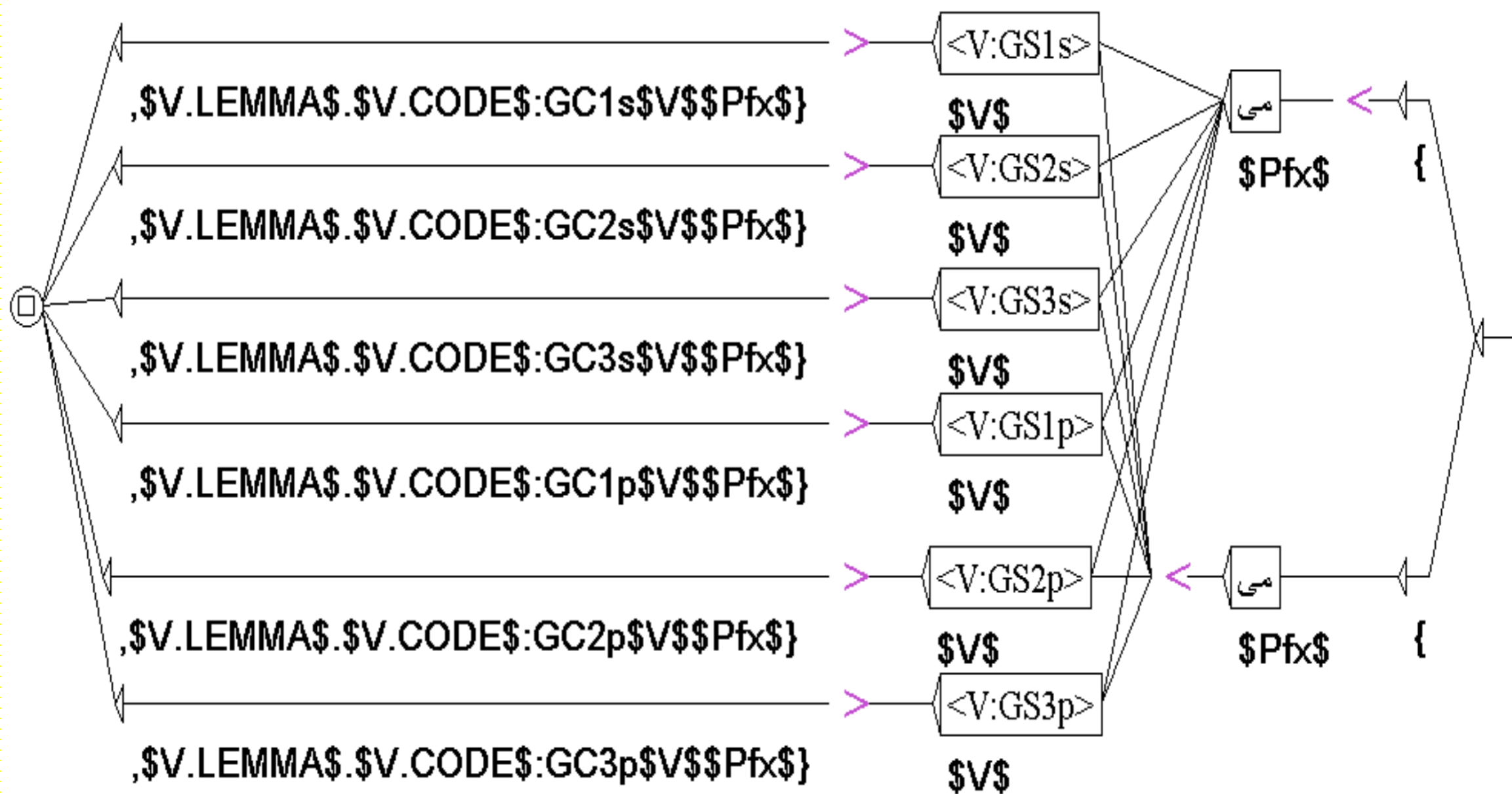


Figure 11. Morphological graph for Past Continuous (GC)

## 4.5.Generating the dictionary of inflected simple forms

<stem>,<inflection-code>+L=<Lemma>

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

رفـتن=L+V002،رفـت
رفـتن=L+V001، رو
آمـدن=L+V002، آمـد
آمـدن=L+V001، آی
شدن=L+V002، شـد
شدن=L+V001، شـو

## 4.6. Building 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.

Each entry in DELAF has the following structure:

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

رفــتم ،رفــت.رفــتن=L+V.رفــت:GS1s
رفــتی ،رفــت.رفــتن=L+V.رفــت:GS2s
رفــت،رفــت.رفــتن=L+V.رفــت:GS3s
رفــتیم ،رفت.رفــتن=L+V.رفــت:GS1p
رفــتیـد،رفت.رفــتن=L+V.رفــت:GS2p
رفــتنـد ،رفت.رفــتن=L+V.رفــت:GS3p
رفــت،رفت.رفــتن=L+V.رفــت:Gr

Figure 12. Extract of Persian DELAF

- 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.

## 4.7. Apply the DELAF to the corpus

The DELAF can now be applied to texts in order to tag the simple verbs. At this stage, 1,536 simple forms of the TEP corpus were tagged by our DELAF, which correspond to 368,831 occurrences and 1,370 different simple word forms (slightly 10% of corpus' simple words).
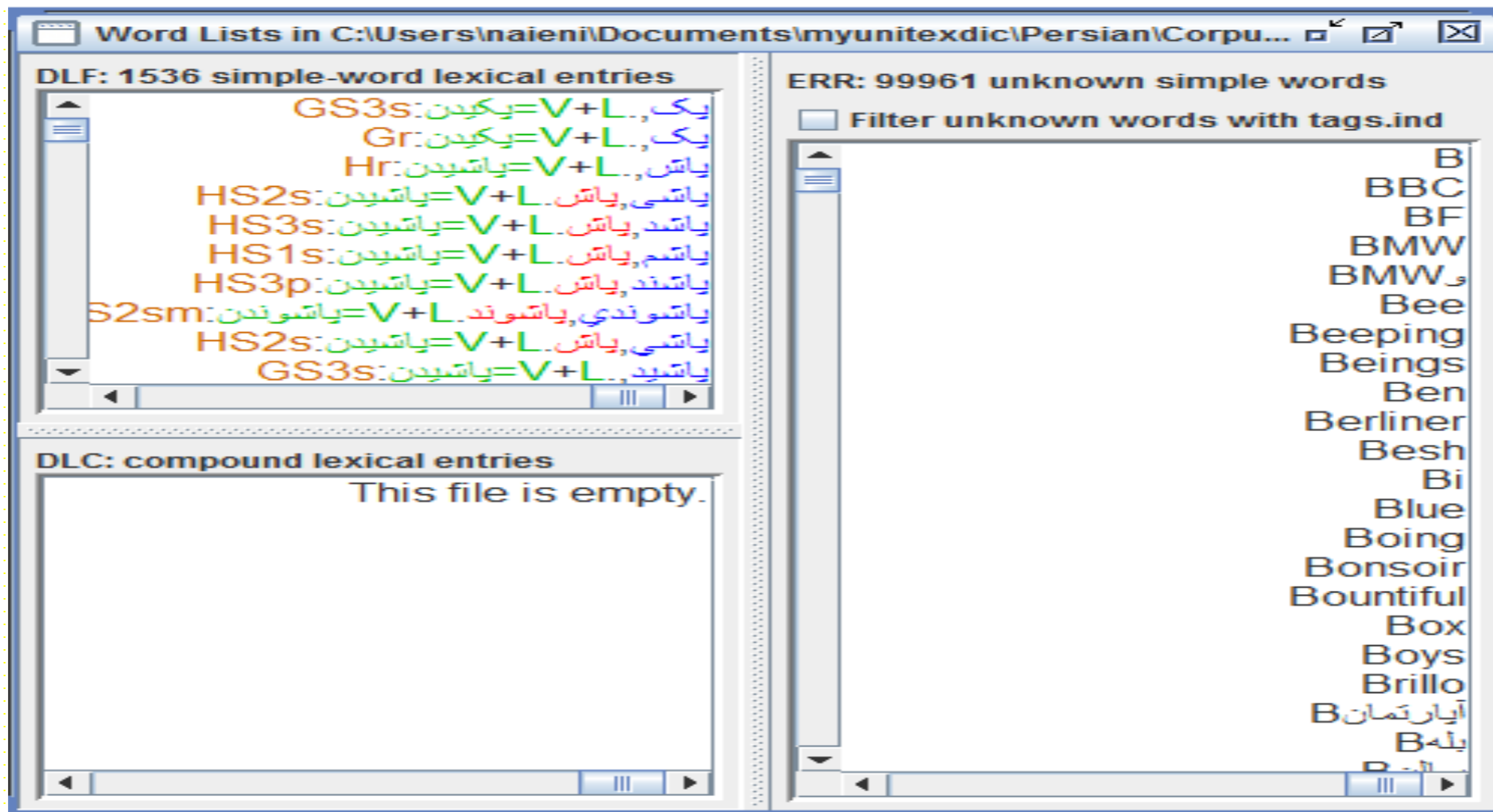


Figure 13. Word list

# 4.8. Apply the compound tenses FST to the corpus

When I applied the compound tenses FST to the corpus, the graph found 82,605 compound tenses matches (1,145 different compound word forms).

| | |
|---|---|
| GP2p :گذاشتن=V+L.کرد,کرده بودید,کرده بودید | HP3s :گذاشتن=V+L.گذاشت, گذاشته |
| GP1p :گذاشتن=V+L.کرد,کرده بودیم | HP3s :گذاشتن=V+L.گذاشت,گذاشته است |
| GP1s :گذاشتن=V+L.کرد,کرده بودم | GP1s :گذاشتن=V+L.گذاشت,گذاشته بودم |
| GP3p :گذاشتن=V+L.کرد,کرده بودند | GP3p :گذاشتن=V+L.گذاشت,گذاشته بودند |
| GF3s :گذاشتن=V+L.کرد, کرده بوده | GF3s :گذاشتن=V+L.گذاشت, گذاشته بوده |
| GF3s :گذاشتن=V+L.کرد,کرده بوده است | HP3s :گذشتن=V+L.گذشت, گذشته |
| HP3s :کندن=V+L.کند, کنده | HP2s :گذشتن=V+L.گذشت,گذشته ای |
| HP3s :کندن=V+L.کند,کنده است | HP3s :گذشتن=V+L.گذشت,گذشته است |
| GP3s :گذاشتن=V+L.گذاشت,گذاشته بود | HP1s :گذشتن=V+L.گذشت,گذشته ام |

Figure 14. first entries of the compound tenses dictionary

# 5.Evaluation

My evaluation was two-fold:

• Evaluate *the formal correction* of the dictionary entries by using the BNFs for both simple and compound forms.

• Estimate *the lexical coverage* of the dictionaries by applying them to a sample text.

## 5.1. Evaluation of the formal correction of the dictionary

The inflected simple forms' dictionary contains 145 lemmas, 292 written verb stems and 127 spoken verb stems and 1,536 different entries. After applying the simple forms BNF (Figure 1) **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, I have 222,698 total entries matches and 3,953 different entries. After applying the Compound form BNF (Figure 2), **all entries were also matched, and no errors were found**.

## 5.2. Evaluation of the lexical coverage of the dictionary

A sample text with around 1,000 words was selected from the web and the text was POS-tagged by the system.

Results:

- For the simple form, we find the verb بودن *bodan* "to be" that appears two times in the text that was not captured by the dictionary due to forgetting to add corresponding inflectional graph.

- In the compound forms' evaluation, after manual verification, we find also the verb نیست *nist* "not to be" that appear one time in the text that was not taken because, we missing to add this compound graph for negation of the verb هست *hast "be"*.

# 6. Conclusion

• I have achieved the main objectives of this project. I built a sufficiently large lexicon of Persian verbs, until I attained a frequency threshold corresponding to a cumulative percentage of 90% of the total corpus word forms (cumulative sum in 3,349,920) corresponding to 1,466 different verb forms.

• I 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.

• Furthermore, I built 22 morphological graphs in order to produce a dictionary of compound verb inflected forms.

- I 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).

- In order to evaluate the correction of the both dictionaries, I applied the simple form BNF and compound form BNF to dictionaries, and no errors were found in both of them.

- In order to evaluate the dictionary lexical coverage, for the simple form, I found the verb بودن *bodan* "to be" that appears two times in the text that was not captured by the dictionary due to forgetting to add corresponding inflectional graph and in the compound form evaluation, after manual verification, I found the verb نیست *nist* "not to be" that appear one time in the text that was not matched because, I was missing to add compound graph for negation of the verb هست *hast "to be"*.

# Future work

Naturally, much is still left to be done.

- First, I 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 the remaining parts of speech (POS), this will be the challenge for future work.

Thank You! :)