# Development of a Persian Syntactic Dependency Treebank

Mohammad Sadegh Rasooli[1], Manouchehr Kouhestani[2], and Amirsaeid Moloodi[3]

[1]rasooli@cs.columbia.edu [2]m.kouhestani@modares.ac.ir [3]a.moloodi@ut.ac.ir

[1]Columbia University [2]Tarbiat Modares University [3]University of Tehran

## Objectives

There was a lack of syntactically annotated data. We tried to create a valuable linguistic data set for the Persian language.

- Second linguistic product by **Dadegan research group** after valency lexicon of Persian verbs [1].
- 30,000 manually annotated sentences.
- The largest syntactic treebank for Persian.
- Extendable to semantic treebank.
- Persian is
  - An Indo-European language.
  - Spoken by more than 100 million speaker.
  - Rich morphology and free word order.
  - *An under-resourced language.*

## Why Dependency Trees?

- Dependency representation is useful for showing
  - Non-projective trees.
  - Compound verbs in Persian.
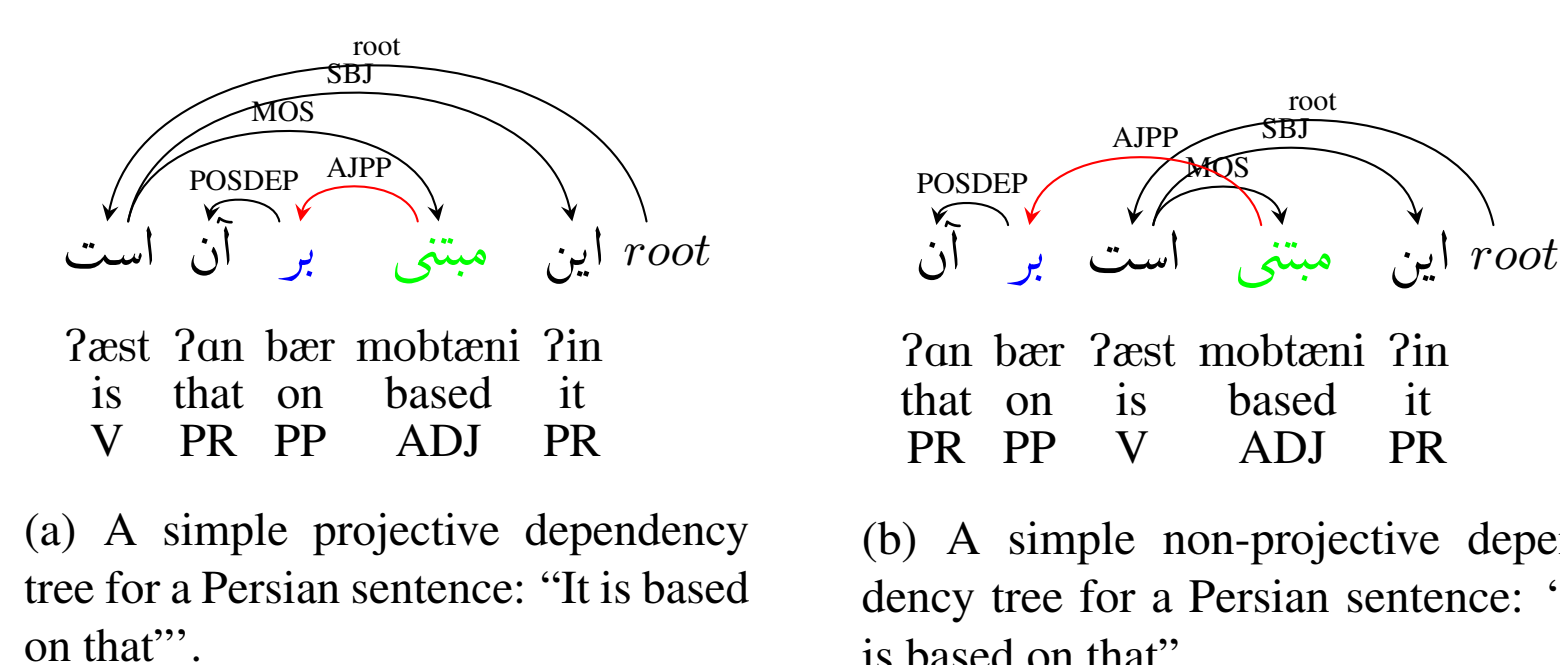- Convertible to phrase-structure trees.



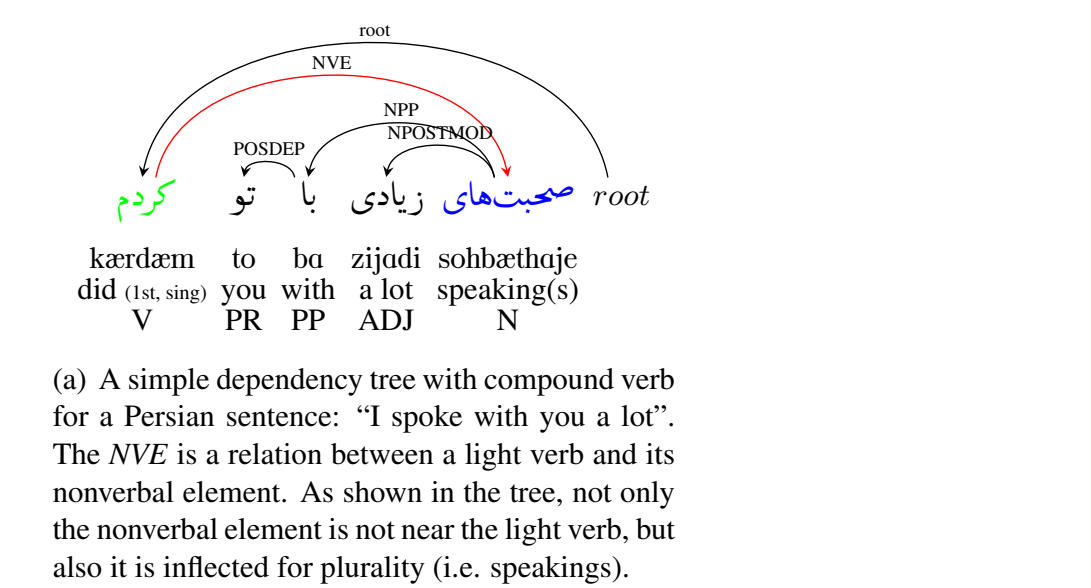Figure 1: An example of free-word order in Persian.



Figure 2: Representation of Persian verbs in dependency trees.

## Annotation Process
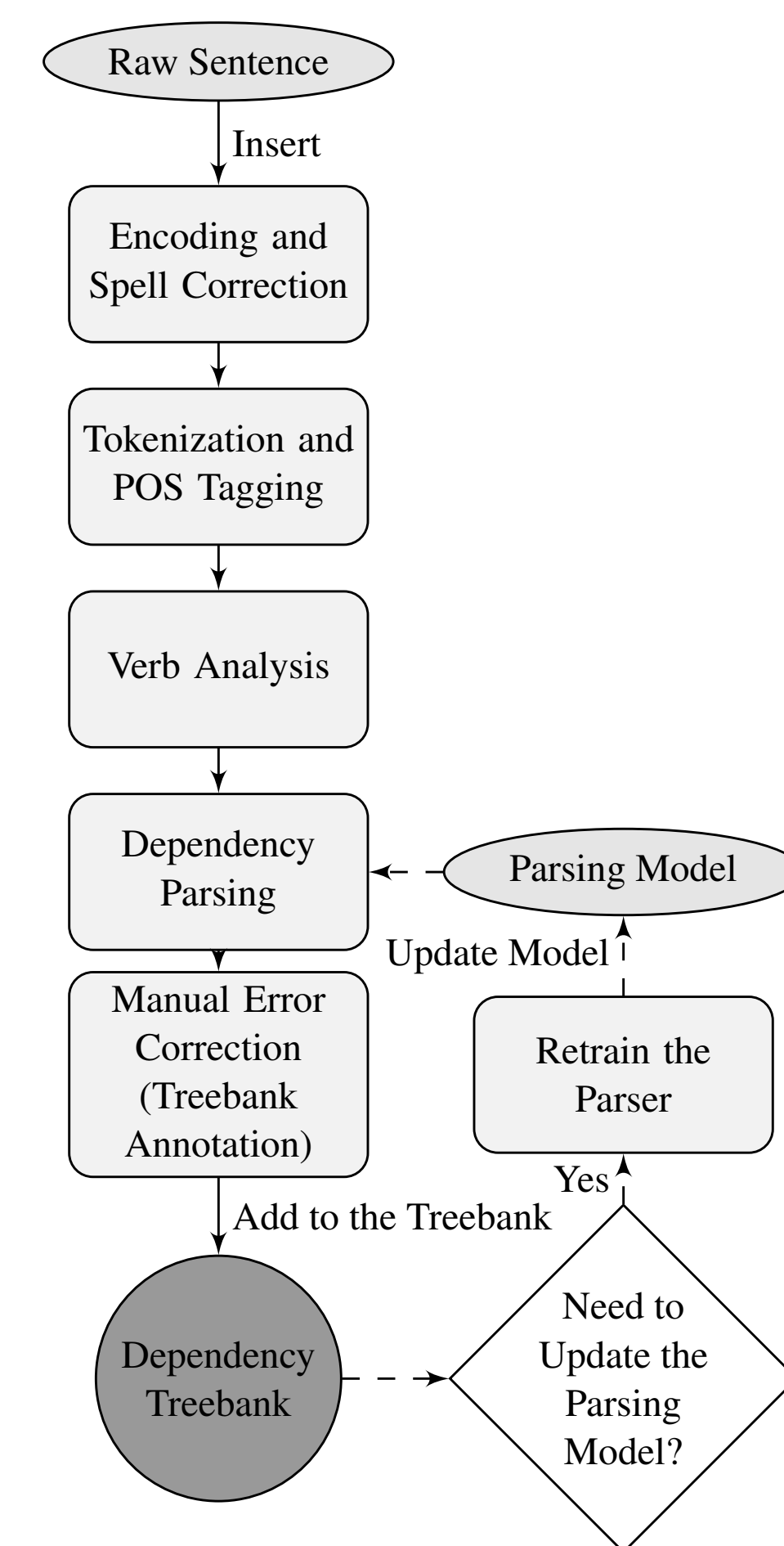
We used a bootstrapping approach to annotate the data.



Figure 3: Diagram of bootstrapping approach in the development of the dependency treebank.

- **Annotators**
  - Started with 8 annotators and gradually increased them to 12 people.
  - 12 months of annotation.
- **Tools**
  - Virastyar [2] for spell checking, lemmatization and POS tagging.
  - Persian verb analyzer [3] for recognizing and lemmatizing verbs.
  - MST parser [4] for parsing input sentences.

## Statistics

- 29,982 sentences; 498K tokens, and 37K types.
- Avg. sentence length: 16.6
- Number of distinct verbs: 4.7K
- 44 dependency relations
- 17 coarse-grained part of speech tags
- 1.8% non-projective sentences (0.02% non-projective arcs)

## Two Different Representations

There are two possible representations for objects accompanied by the case marker:

- Case marker as a post-position is the head of the object phrase.
  - Creates more non-projective trees.
  - Simplifies the search for objects (closer to the verb than the object and the object should come before it).
- Object is the head of the case marker.
  - Closer to the human interpretation.
  - This representation is provided by automatic conversion from the first representation
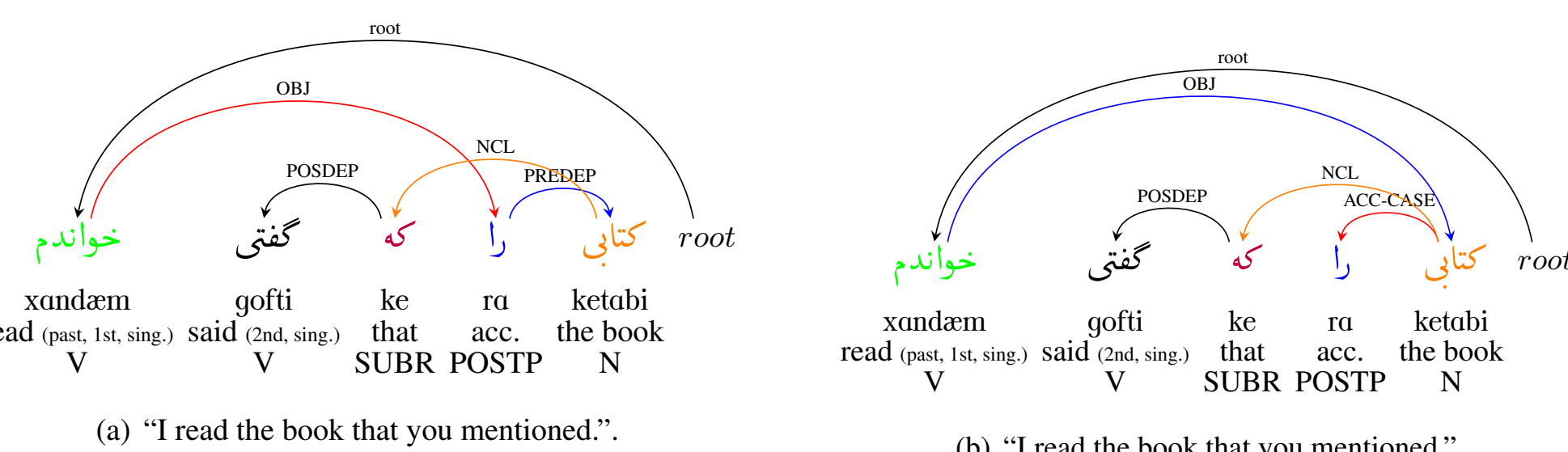


Figure 4: A sample sentence with two kinds of representations of object-verb relation.

## Correcting Potential Errors

- We provided additional scripts for finding potential errors;
  - E.g. annotation mistakes such as a verb as a subject for a noun or mismatch with valency lexicon information.

| | |
|---|---|
| Changes to Unlabeled Relations | 4.91% |
| Changes to Labeled Relations | 6.29% |
| Changes to POS Tags | 4.23% |

Figure 5: Statistics about changes in the treebank after the manual correction of the potential errors.

## Annotators' Agreement

5% of the data is doubly annotated.

| | |
|---|---|
| Unlabeled Relations | 97.06% |
| Labeled Relations | 95.32% |
| POS Tags | 98.93% |

Figure 6: Statistics about agreements among the annotators.

## Future Direction

- Create other resources such as SRL treebank.

### Online Treebank Search

An online tool for searching dependency relations



Figure 7: Dadegan dependency treebank search tool.

http://search.dadegan.ir/advance/

## References

[1] M. S. Rasooli, A. Moloodi, M. Kouhestani, and B. Minaei-Bidgoli. A syntactic valency lexicon for Persian verbs: The first steps towards Persian dependency treebank. In *LTC*, pages 227–231, Poznań, Poland, 2011.

[2] O. Kashefi, M. Nasri, and K. Kanani. *Automatic Spell Checking in Persian Language*. SCICT, 2010.

[3] M. S. Rasooli, H. Faili, and B. Minaei-Bidgoli. Unsupervised identification of Persian compound verbs. In *MICAI*, pages 394–406, Puebla, Mexico, 2011.

[4] R. McDonald, K. Crammer, and F. Pereira. Online large-margin training of dependency parsers. In *ACL*, pages 91–98, Sydney, Australia, 2005.

## Contact Information

- Web: http://www.dadegan.ir/en/perdt
- Email: info@dadegan.ir

## Acknowledgements