

Language Learning from Commentaries in Videos using Dynamic NLP

Anusha Chowdhury, Kaviti Sai Saurab
Guide: Prof. Amitabha Mukherjee

Department of Computer Science and Engineering, Indian Institute of Technology Kanpur

Problem

Learning a new language is tough even for a human being. Here we explore different techniques of Dynamic NLP through which an intelligent agent can learn a new language like Bengali or Telugu (word, pattern and syntax learning) from video commentaries.

Our Approach

- We collected 18 commentaries in Bengali and 9 in Telugu for a particular video.
- Next step was to construct a histogram for each data set and group them into clusters by applying the **Bhattacharyya distance**.
- The corpus is now divided into:
 - Family-lect** : commentaries which have a lot of similarity
 - Multi-lect** : the variations come under this category.
- Now perform **Contrastive Association** and system learns some high-confidence words for each concept.
- To **prune** the corpus, we cut down the words with high frequency and low contrastive score or extremely low freq.
- We use **word2Vec model** to learn the synonyms present in the corpus.
- For pattern and syntax learning, we applied the **ADIOS** algorithm.

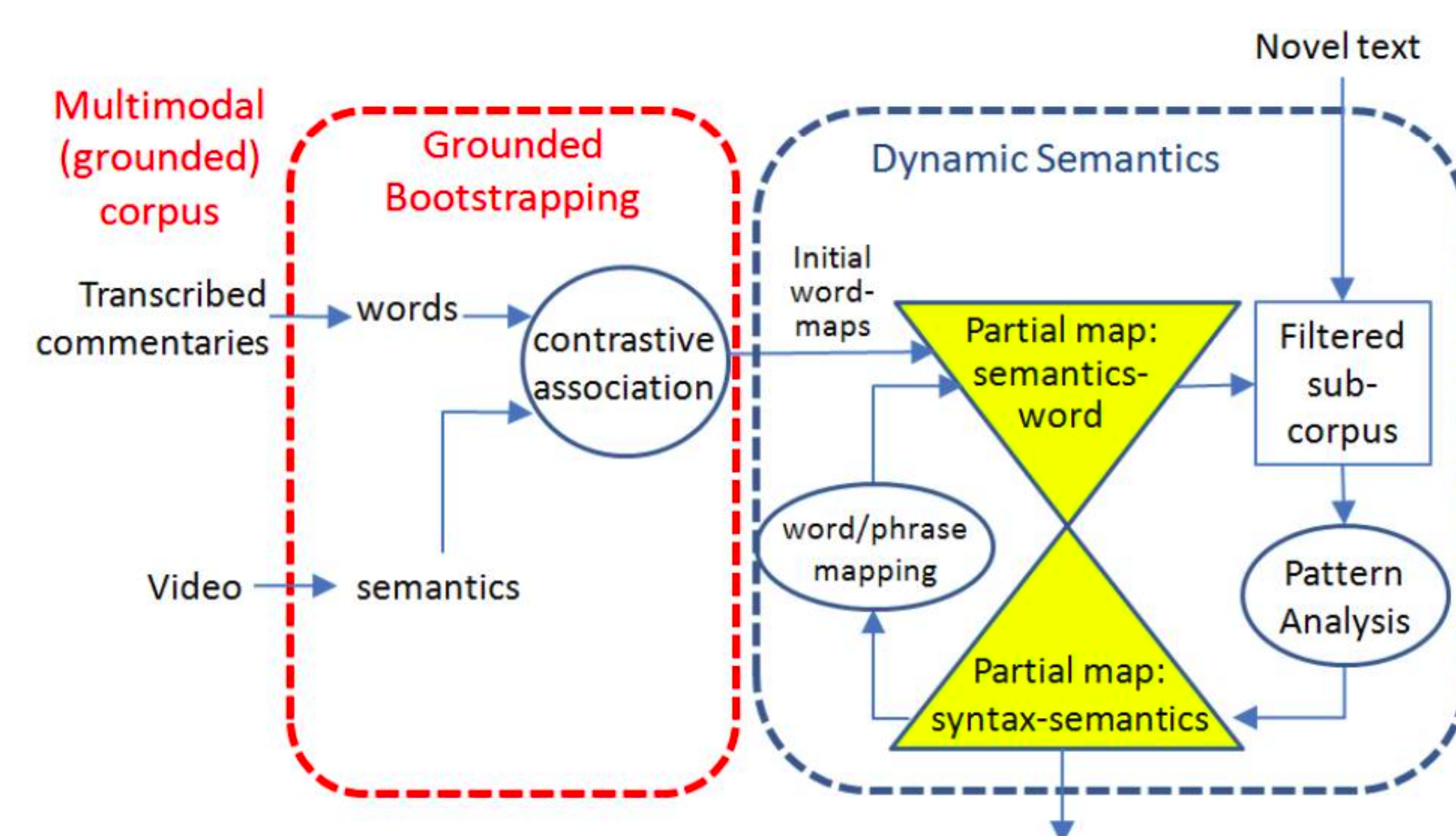


Figure 1: Dynamic NLP approach overview [1]

Contrastive Association

For disjoint concepts c_1, c_2 we divide the sentences into sets - those that occur in commentaries for video involving c_1 , and those that occur for c_2 [4].

$$Score_{w_i, c_1} = P(w_i, c_1) / P(w_i, c_2) \quad (1)$$

The following results were obtained for Bengali:

Contrastive Association (only top 2 or 3 shown):			
Schema	Top Bengali Lexeme(Freq/C contrastive Score)	OppSchema	Top Bengali Lexeme(Freq/Pr obability Score)
AGENT(Dome)	ডোম (27) ডোম (3) (3.0) পৌষাণ (9)(2.25)	AGENT(Daisy)	ডেই (21) 0.11 বল (14) 0.07 প্রকটা (13) 0.06
OBJECT(Ball)	বল (27) অবধি (5)(5.0) না (10)(3.33)	OBJECT(Square)	বাক্স (14) 0.08 প্রকটা (14) 0.07 লাল (12) 0.06
COLOUR(Red)	লাল (24) ফিল (2)(2.0) লাক্সার (2)(2.0)	COLOUR(Blue)	নীল (24) 0.11 বল (15) 0.07 প্রকটা (14) 0.07
ACTION(Throw)	ফুটবল (21) লাক্সার (16)(3.2)	ACTION(Roll)	গড়িয়ে (14) 0.08 ফিল (13) 0.07
PATH(Before target)	না (11) ফিল (10)(10.0)	PATH(On/After target)	প্রকটা (24) 0.07 গড়িয়ে (12) 0.03

Figure 2: Contrastive Association on Bengali, Telugu

Histogram Similarity using Bhattacharyya distance

To measure the similarity of two discrete probability distributions, we calculate the Bhattacharyya coefficient as:

$$Bhattacoeff = \sum_{i=1}^n \sqrt{(\sum a_i)(\sum b_i)} \quad (2)$$

where samples a and b have n partitions, and $\sum a_i, \sum b_i$ are the no. of members of samples a and b in i^{th} partition. For discrete distribution $Bhattadistance = -\ln(Bhattacoeff)$.

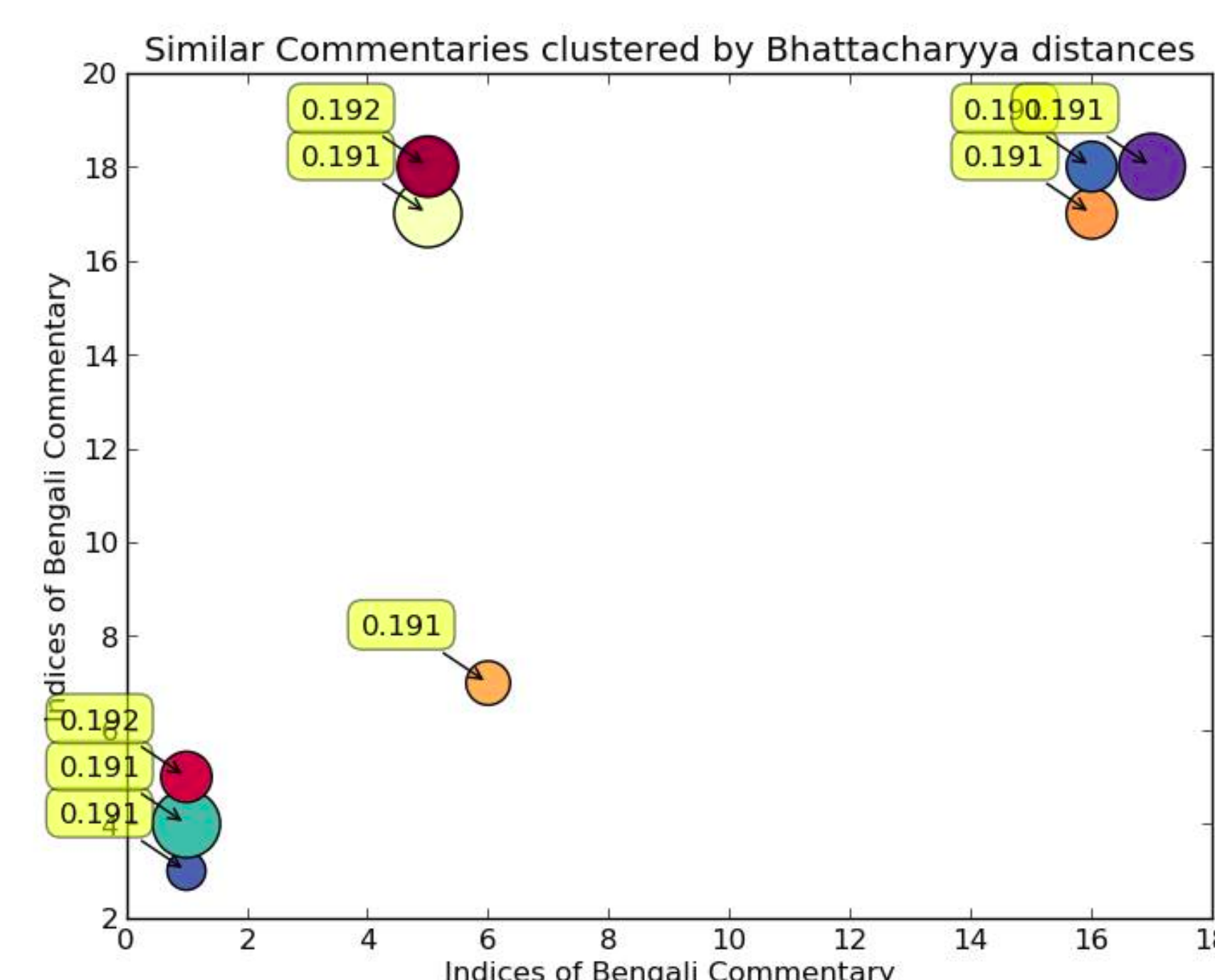


Figure 3: Clustering of Histograms from Bhattacharyya dist.

ADIOS(Automatic Distillation of Structure)

Adios works in an unsupervised fashion to extract significant patterns from a raw unannotated corpus. It works based on two principal components.

★ Representational Data Structure

★ Pattern Acquisition Algorithm

Probabilistic inference of pattern acquisition is used to cluster similar lexicons and recursive construction is used to hierarchically generate more complex patterns. Significance of patterns is determined using the **motif extraction criterion** [3]. At the end of each iteration, the most significant pattern is added to the lexicon as a new unit and the graph is rewired by merging the paths it subsumes. Algorithm stops when new patterns or equivalence classes are not discovered anymore.

Bengali		Telugu	
নীল	লাল	ఎర్ర	నీలం
[blue]	[red]	[red]	[blue]
বল	বাক্স	దగ్గరలో	দূরারলো
[ball]	[box]	[near]	[far]
অবধি	পর্যন্ত	చదరవు	ఒక్కొక్కటి
[till]	[till]	[square]	[box]

Figure 4: Equivalence classes obtained from ADIOS

Applying word2vec model on our corpus

The word2vec tool provides an efficient implementation of continuous **bag-of-words** and **skip-gram** architectures for computing vector representations of words. Here we apply it on our small corpus and interestingly we get quite good results in Bengali. It has learnt synonyms and also it can find the odd word out from a given list.

```
model.similarity("লাল".encode('utf-8'), "নীল".encode('utf-8'))
returns 0.94 whereas for dissimilar words it is around 0.25
model.doesnt_match("লাল না নীল সবুজ".encode('utf-8').split())
returns না
```

Figure 5: Results from word2vec

RESULTS OF ADIOS

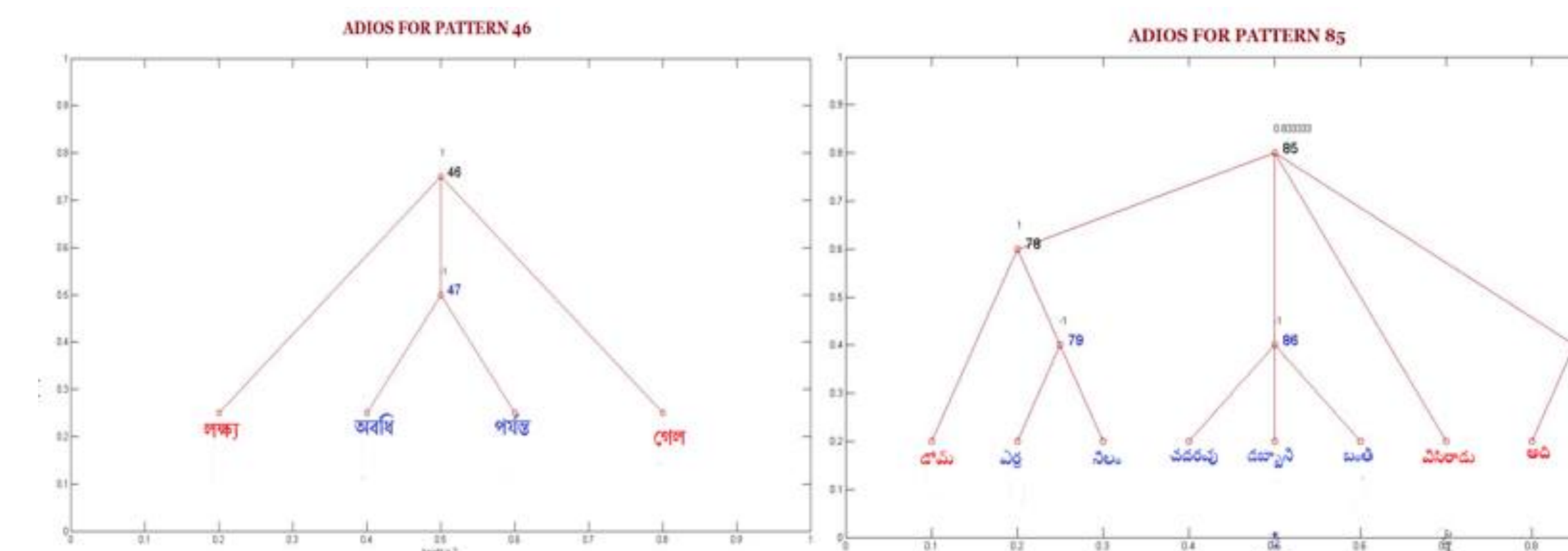


Figure 6: Results from ADIOS for Bengali, Telugu

Conclusion

Comparison with already existing work:

- [1] applies dynamic NLP on Hindi and English languages. Both [1] and [2], the grounding of target did not achieve high confidence. But Telugu, Bengali corpus yielded high confidence association of on, after target.
- Apart from ADIOS we also used word2Vec to study equivalence classes and performance was really good.
- We are implementing **Morphosyntactic discovery** as final step where we will do a text-based morphological similarity analysis using Levenshtein distances between words (for example, in Bengali *lokhyo* and *lokhye* both refer to the same semantic concept target). These variations in pattern can be due to tense or gender difference in a language.

References

- [1] D. Semwal. Dynamic NLP : A model for language bootstrapping and lifelong learning M.Tech. Thesis under Dr. A. Mukerjee at IIT Kanpur, 2014.
- [2] D. Semwal, S. Gupta, A. Mukerjee. From visuo-motor to language.(2014)
- [3] Z. Solan, D. Horn, E. Ruppim, S. Edelman. Unsupervised learning of natural languages(2005).
- [4] E. M Markman. Constraints children place on word meanings.