

**CS365 Project**  
**LANGUAGE LEARNING FROM BENGALI**  
**COMMENTARY IN VIDEOS**

**Guide: Prof. Amitabha Mukherjee**  
Anusha Chowdhury

March 28, 2015



# Problem Statement

Introduction

**Problem Statement**

Our Approach

Related Works

CONCLUSION

The system takes a set of commentaries on videos as input for the purpose of **Bengali language learning**.

The main aim is **word learning** and **syntax learning**. Syntax basically refers to the positioning of subject, object and verb. The input videos have agents performing some actions on coloured objects and there is a target.

# Motivation

Main motivation was the urge to learn Bengali language from scratch just as a child does.



**Figure :** A child learns language by **semantic mapping**: how linguistic elements relate to visible situations

# DYNAMIC NLP: Preliminaries

Introduction

**Problem  
Statement**

Our Approach

Related Works

CONCLUSION

- Process of learning involves continuous expansion of the already acquired language model.
- There is continuous learning from new sentences
- All evaluations of sentence semantics is partial
- Associations between word or phrase and meaning is dynamic which may broaden or shrink.

# Our Approach

Introduction

Problem  
Statement

**Our Approach**

Related Works

CONCLUSION

- Collect the input data (Bengali commentary)
- Identify the **Family-lect** and the **Multi-lect** from the corpus
- Perform **contrastive association** based on probability calculations
- Words like 'and' are ignored to shorten down the corpus
- Patterns are learnt by applying **ADIOS** algorithm on the family lect after pruning the corpus
- Learning the sentence syntax from verb phrases
- **Morphosyntactic discovery**: Text-based morphological similarity analysis gives rise to clusters.

# The video

A glimpse of the video, each having 16 frames:

Introduction

Problem  
Statement

**Our Approach**

Related Works

CONCLUSION

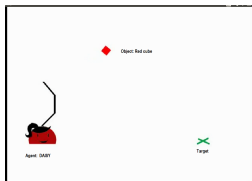


Figure : Daisy is throwing a Red Cube

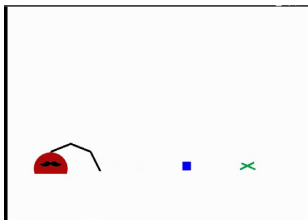


Figure : Dome is rolling a blue cube

# Pre-linguistic concepts

Introduction

Problem  
Statement

**Our Approach**

Related Works

CONCLUSION

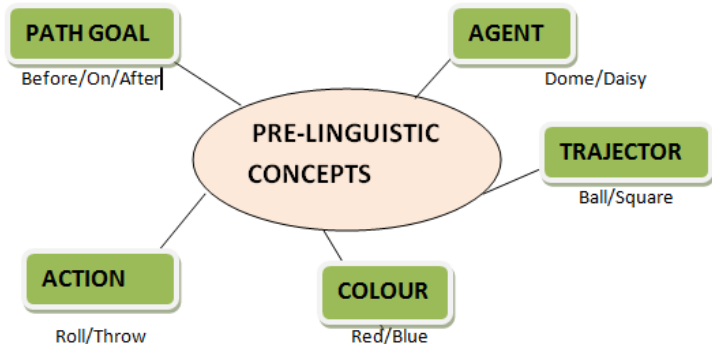


Figure : The pre-linguistic concepts

# The interface

Introduction

Problem  
Statement

Our Approach

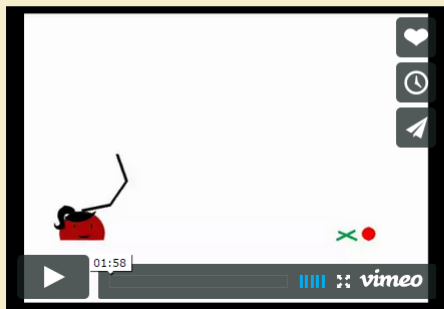
Related Works

CONCLUSION

A snapshot of the interface for recording commentary in Bengali:

Press the button below to start recording your commentary

Start Video and Recording





# The commentaries

A sample of the Bengali commentaries collected:

FRAME NO.	Text in Bengali
1 67-247	ডোম লাল বাক্স ছুঁড়ল কিন্তু সেটা লক্ষ্যে পৌঁছালো না
2 248-455	ডেজী লাল বাক্স গড়িয়ে দিল, লক্ষ্যের আগে চলে গেল
3 456-656	ডোম নীল বল ছুঁড়ল কিন্তু লক্ষ্যে পৌঁছালো না
4 657-821	ডেজী লাল বাক্স ছুঁড়ল লক্ষ্যে পেরিয়ে চলে গেল
5 822-1005	ডোম নীল বল গড়িয়ে দিল লক্ষ্যে পেরিয়ে গেল
6 1006-1168	ডোম নীল বাক্স গড়িয়ে দিল ঠিক লক্ষ্যে পৌঁছালো
7 1169-1339	ডোম লাল বল গড়িয়ে দিল ঠিক লক্ষ্যে গেল
8 1340-1548	ডেজী নীল বল গড়িয়ে দিল আবার ঠিক লক্ষ্যে পৌঁছালো
9 1549-1684	ডেজী লাল বল ছুঁড়ে দিল লক্ষ্যে পেরিয়ে গেল
10 1685-1878	ডোম লাল বাক্স ঠেলার চেষ্টা করল
11 1879-2066	ডোম নীল বল ছুঁড়ল ঠিক লক্ষ্যে গেল
12 2067-2214	ডেজী লাল বল ছুঁড়ল লক্ষ্যে অবধি গেল না
13 2215-2410	ডোম লাল বল ছুঁড়ল লক্ষ্যে অবধি গেল না
14 2411-2588	ডেজী একটা নীল বাক্স ছুঁড়ল, ঠিক লক্ষ্যে পৌঁছালো
15 2589-2783	ডোম নীল বাক্স ছুঁড়ল, লক্ষ্যে পেরিয়ে চলে গেল
16 2784-2947	ডেজী নীল বল ছুঁড়ল কিন্তু লক্ষ্যে অবধি গেল না

# Family-lect and Multi-lect

Introduction

Problem  
Statement

**Our Approach**

Related Works

CONCLUSION



- Family-lect : the set of commentaries which are coherent with respect to the lexical choices for trajectors, actions, agents and the constructional choices.
- Multi-lect: this corpus includes the different varieties of syntax and vocabulary available.

# Contrastive Association

Introduction

Problem  
Statement

Our Approach

Related Works

CONCLUSION

For two non-overlapping concepts  $c_1, c_2$  we first divide the corpus into subsets:

- 1 those that arise on commentaries for video involving  $c_1$
- 2 those arising for  $c_2$

Now a scoring function for association is defined as the ratio of the joint probabilities of word  $w$  occurring with concept  $c_1$  and that with  $c_2$ :

$$S_{w,c_1} = \frac{P(w, c_1)}{P(w, c_2)}$$

An overview of ADIOS(Automatic Distillation Of Structure).  
Algorithm:

- Initialization :loading all sentences
- for all  $i = 1 : N$   
    Pattern Distillation( $i$ )  
    Generalization( $i$ )  
    endfor
- repeat until no more significant patterns are found

# Related Works

Introduction

Problem  
Statement

Our Approach

**Related Works**

CONCLUSION

- Mukherjee et. al.[From visuo-motor to language, 2014] have shown how how a learning agent learns syntactic patterns based on some highconfidence words for English and Hindi.
- In D.Semwal's thesis [Dynamic NLP, 2014] we get an overview of the different dynamic NLP techniques that can be used.
- Z.Solan et. al. have outlined the ADIOS algorithm in[14].
- Zettlemoyer et. al. delve deep into morphosyntactic discovery in[14].

# References

Introduction

Problem  
Statement

Our Approach

Related Works

CONCLUSION

- D. Semwal, S. Gupta, A. Mukerjee. From visuo-motor to language (2014) At [http://www.cse.iitk.ac.in/users/amit/pub/aaai-fs\\_14.pdf](http://www.cse.iitk.ac.in/users/amit/pub/aaai-fs_14.pdf).
- D. Semwal. Dynamic NLP : A model for language bootstrapping and lifelong learning M.Tech. Thesis under Dr. Amitabha Mukerjee (2014). At <http://www.cse.iitk.ac.in/users/deepalis/thesis/report.pdf>.
- Z. Solan, D. Horn, E. Ruppin, S. Edelman. Unsupervised learning of natural languages. Proceedings of the National Academy of Sciences pp:11629-11634, 2005.
- A. Wang, T. Kwiatkowski, L. Zettlemoyer. Morpho-syntactic Lexical Generalization for CCG Semantic Parsing. In EMNLP, 2014.

# CONCLUSION

Introduction

Problem  
Statement

Our Approach

Related Works

**CONCLUSION**



# QUESTIONS

Introduction

Problem  
Statement

Our Approach

Related Works

**CONCLUSION**

