

# Antarlekhaka: A Comprehensive Tool for Multi-task Natural Language Annotation

# Hrishikesh Terdalkar

# **Arnab Bhattacharya**

hrishirt@cse.iitk.ac.in

Yoda, Star Wars

arnabb@cse.iitk.ac.in

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

#### Aim

One-stop solution for Natural Language Annotation

#### Motivation

- World's ~7000 languages are low-resource languages
- Human annotation remains extremely relevant

#### **Linguistic Phenomena**

- Ambiguous or absent sentence boundaries
- Rearranging, splitting or merging tokens may be required
- Limited support in existing annotation tools

#### **Punctuation and Word Order**

- " if no mistake you have made losing you are a different game you should play "
- 'If no mistake you have made, losing you are. A different game you should play. "
- "If you have made no mistake, you are losing. You should play a different game. "

#### Sanskrit Example

- Majority of Sanskrit literature in poetry format
- Sentence boundary and token manipulation required

[ na rocate mama-api-etad-ārye ] [ yad-rāghavo vanam / tyaktvā rājyaśriyam gacchet ] [ striyā vākyavaśam gataḥ // 2 viparītas ca vṛddhas ca viṣayais ca pradharṣitaḥ / nṛpaḥ kim iva na brūyāc codyamānaḥ samanmathaḥ // 3 ] 3 L ... J

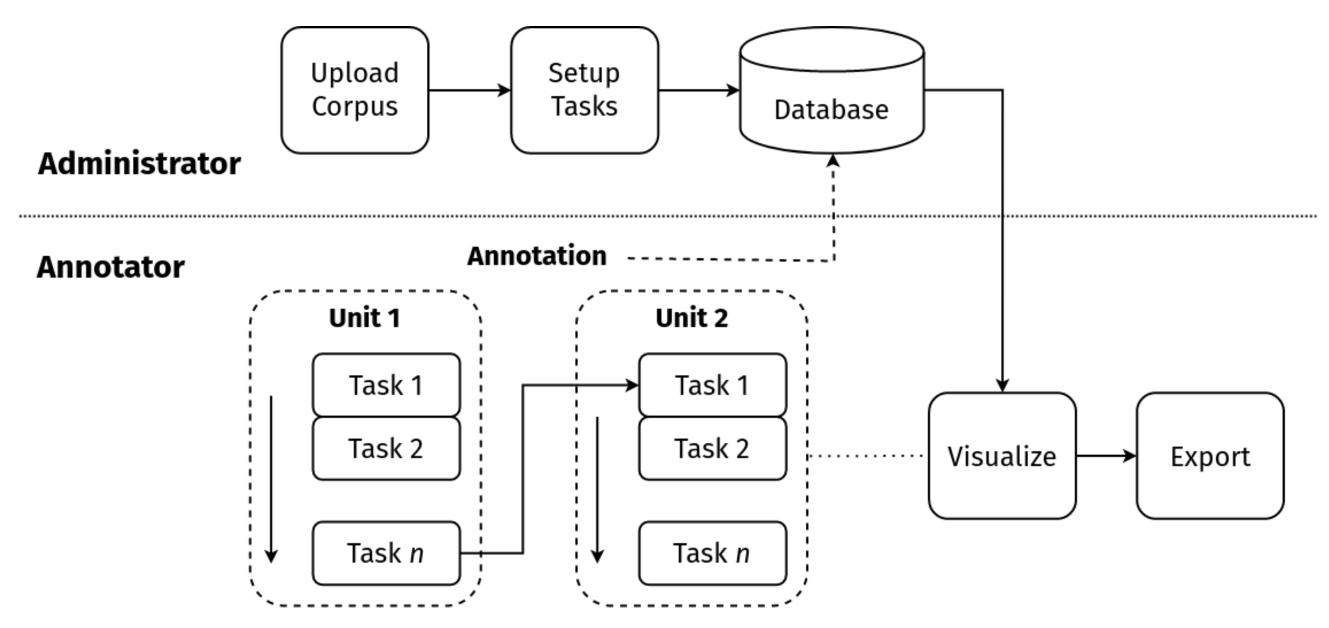
[ ārye etad mama api na rocate ] 1 [ yad rāghavo rājyaśriyam tyaktvā vanam gacchet ]2 viparītah vṛddhah ca viṣayaih pradharṣitah ca codyamānah samanmathaḥ ca striyā vākyavaśaṃ gataḥ nṛpaḥ kim iva na brūyāt ]3 [ ... ]

**Example from Rāmāyaņa** 

### Requirements

- Intrinsic support for sentence boundaries and token order
- Comprehensive task coverage and task customization
- Multiple annotation tasks for same text

# **Architecture**



**Workflow of Antarlekhaka** 

#### **Features**

- Sequential annotation towards multiple NLP tasks
- **Eight categories** of tasks  $\implies$  Task-specific annotation interfaces
- Pluggable heuristics to aid annotators
- Task Management, Ontology Management, Progress Report, Clone Annotations
- Export in Human-readable and Machine-readable format
- Language agnostic, Unicode support

#### **Task Categories**

#### **Sentence Boundary Detection**

- Languages such as Sanskrit
- Corpora in poetry format

#### Token Manipulation: Addition, Exclusion, Merging, Splitting, Ordering

- Canonical Token Order
- Word Segmentation
- Word Grouping

#### **Token Annotation Token Classification** Named Entity

- Lemmatization
- Morphological Analysis
- Spelling Correction
- Part-of-speech Tagging Phonetic Transcription Compound Classification

#### **Token Connection**

- Co-reference Resolution
- Interaction Networks Text Clustering
- **Sentence Classification**
- Sentiment Detection

Recognition

 Sarcasm Detection Spam Detection

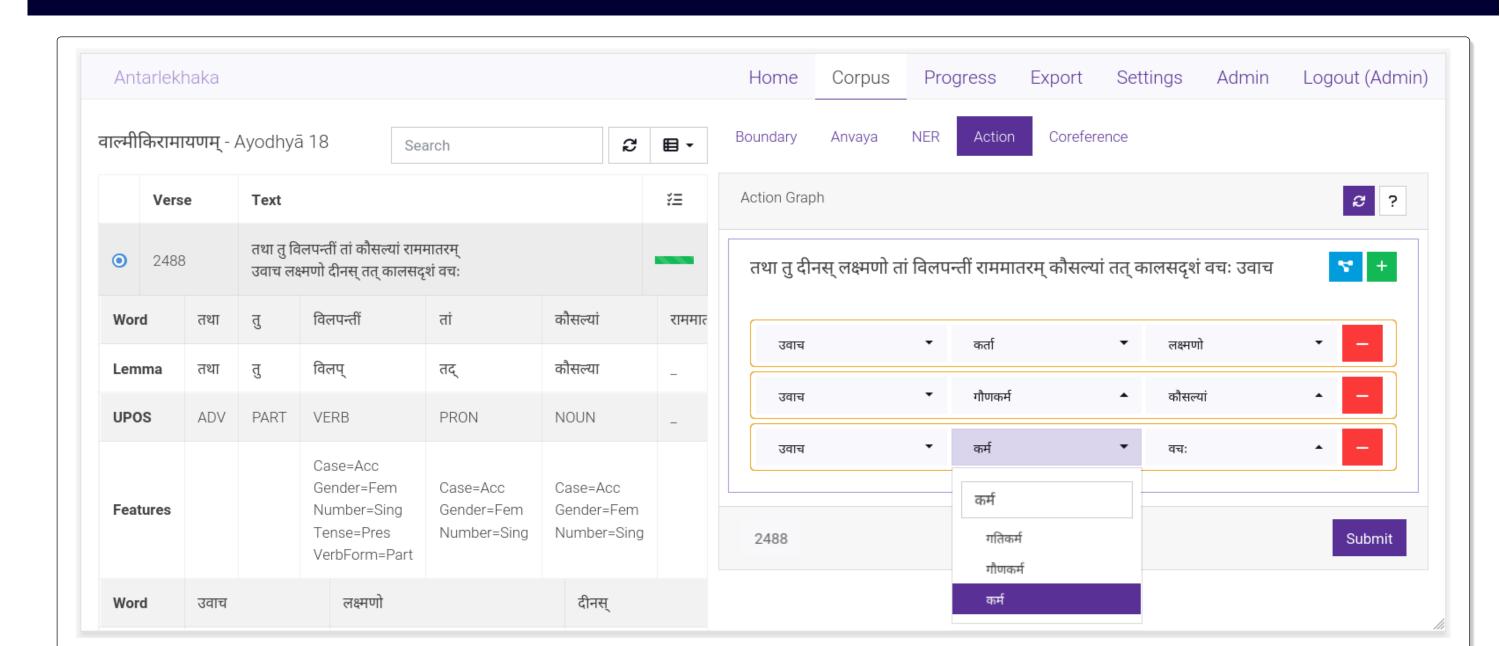
# **Token Graph**

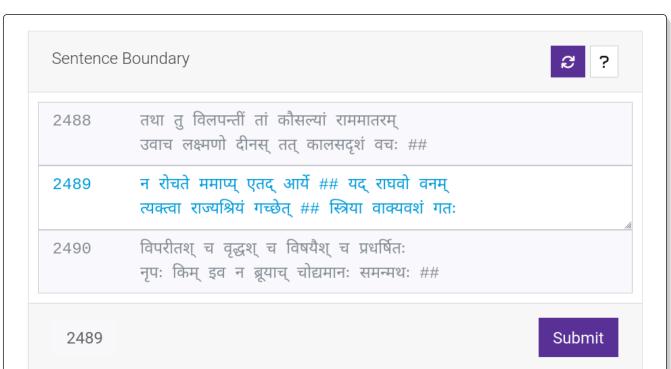
- Dependency Parsing
- Constituency Parsing
- Semantic Graph Action Graph

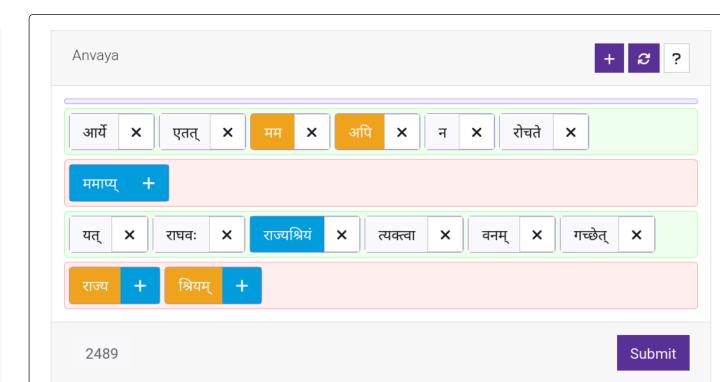
#### **Sentence Graph**

- Discourse Graph
- Timeline Annotation

# **Annotation Interface**

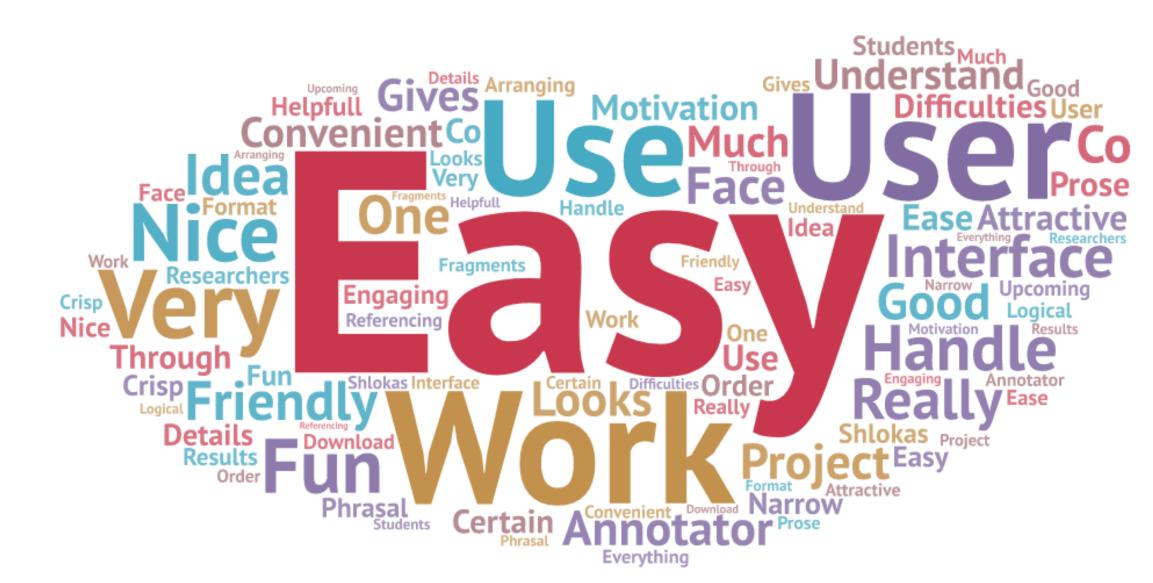






# **Evaluation**

- Total 29 Criteria: Technical, Functional, Data Related, Task Related
- Scores: Antarlekhaka (o.79), INCEPTION (o.74), Sangrahaka (o.74), FLAT (o.71)
- Only Antarlekhaka supports token ordering



**Wordcloud of Survey Responses** 

#### **Open Source Software**

Antarlekhaka/code



#### Acknowledgements

We thank Chaitali Dangarikar, Shubhangi Agarwal, V S D S Mahesh Akavarapu, and Pralay Manna for their valuable feedback.















