# End-to-End Natural Language Understanding Pipeline for Bangla Conversational Agents

**Fahim Shahriar Khan[1] , Mueeze Al Mushabbir[1] , Mohammad Sabik Irbaz[2] , MD Abdullah Al Nasim[2]**

Department of Computer Science and Engineering, Islamic University of Technology[1], Machine Learning Team, Pioneer Alpha Ltd.[2]
khanfahimshahriar0@gmail.com, almushabbir@iut-dhaka.edu, sabikirbaz@iut-dhaka.edu, nasim.abdullah@ieee.org
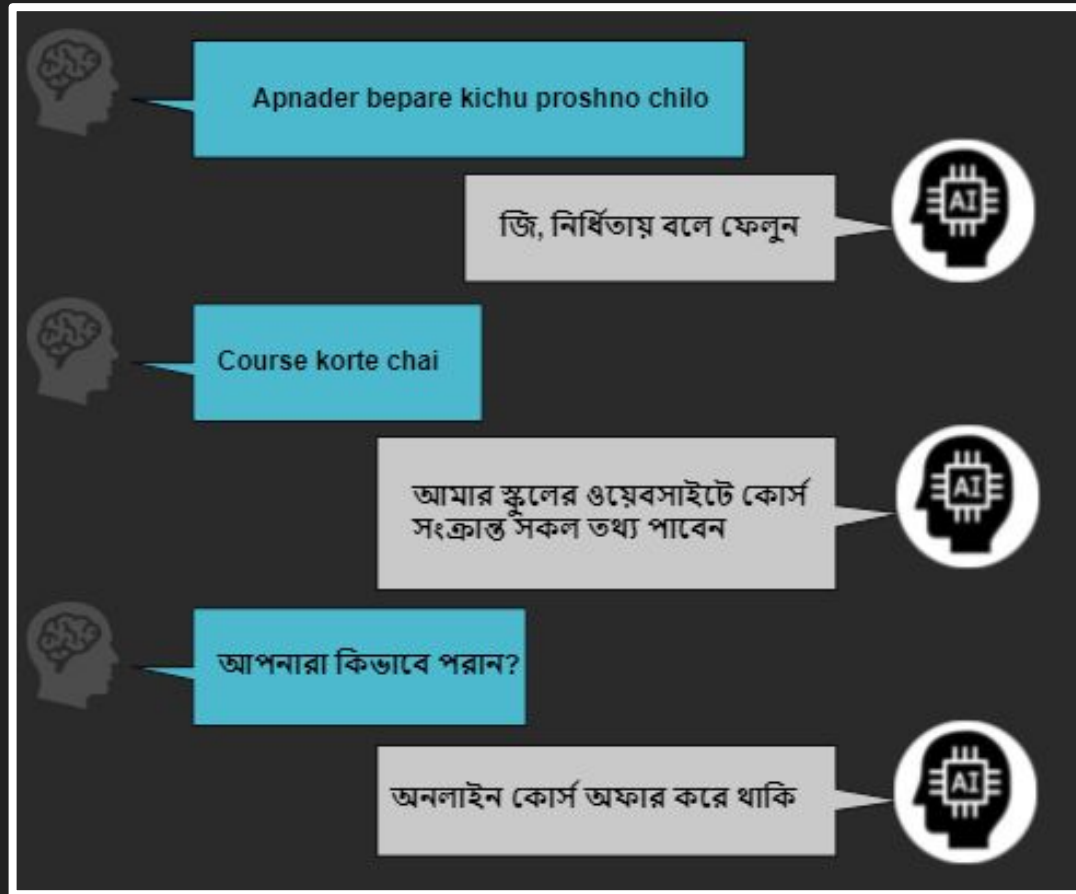
Paper ID : 341

# Introduction

- **Conversational AI Agents aim to provide virtual assistant like services in the form of a dialog system using natural languages.**

- **Existing chatbot systems usually do not have enough support for low-resource languages, like Bangla.**

- **We aim to build an end-to-end (from creating a dataset to evaluating components and pipelines) Bangla Chatbot - to be used as a virtual assistant in the business environments.**

# Example

# Motivation

Support for existing dialog systems in low-resources languages like Bangla, or Bangla Transliteration is scarce.

**+**

With the increasing need of customer services in the virtual world, businesses look for automated service providing solutions.

**+**

Comparative analysis among various components and pipelines with proper reasoning are difficult to figure out.

# Problem Statement

Using low-resource, scarce and skewed Bangla dataset for intent recognition and entity extraction, we need to build an end-to-end NLU pipeline for chatbots which can receive messages and send responses seamlessly as a Business Assistant.

# Research Challenges

**Skewed Low-Resource Data**

1

Available datasets for languages like Bangla are low-resource, low-quality skewed datasets

**Dealing with Bangla Transliteration in English**

2

[describe transliteration]

**Technical Analysis of Components and Pipelines**

3

Understanding the technical properties of each component and pipeline, and their comparative performances with proper reasonings are difficult to figure out.
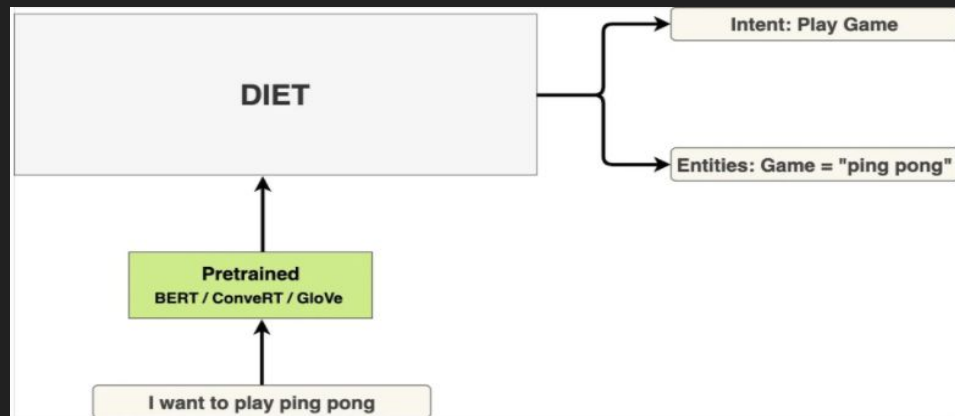
# Related works

## RASA

# Related works

## DIET Classifier

- Dual Intent Entity Transformer (DIET)
- Multi-task transformer architecture
- used for both _Intent Classification_ and _Entity Extraction_
- Provides Modularity
  - used with various pre-trained embeddings like BERT, GloVe, etc.
- Comparable Performance against large-scale pre-trained language models
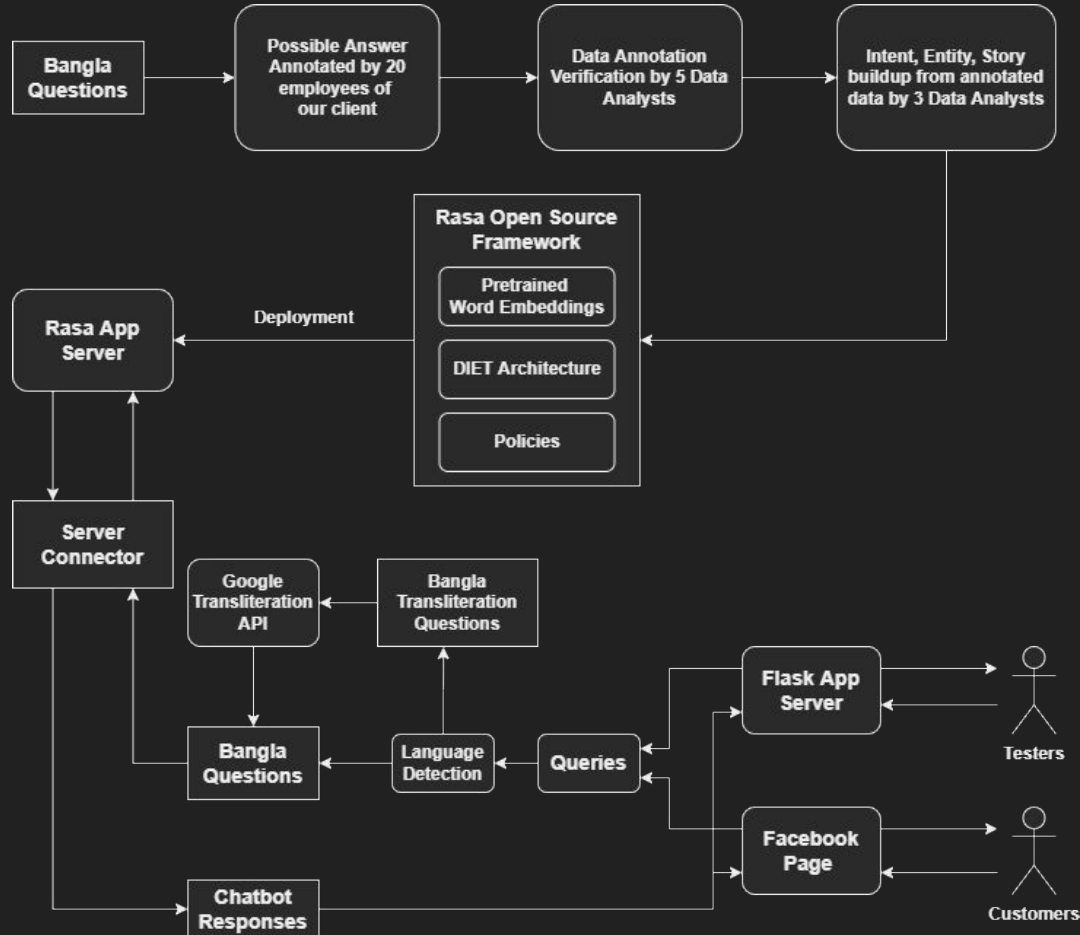  - Faster and Better than BERT



https://rasa.com/blog/introducing-dual-intent-and-entity-transformer-diet-state-of-the-art-performance-on-a-lightweight-architecture/

# Related works

[Vietnamese Paper]

# Machine Learning Life Cycle

# Dataset Preparation

# Dataset overview

- Rasa open source architecture expects the dataset to be partitioned into a 3 separate yml files.
  - nlu.yml
  - domain.yml
  - stories.yml

- Our chatbot needs to deal with FAQs in Bangla and Bangla Transliterations in English.

- FAQs were collected from our client's interaction history with customers from different social media platforms.

- They were annotated by 20 employees of our client.

# Parts of the Dataset

1. *nlu.yml*
   - The FAQs gathered from our client is labelled into intents and entities.

   - In our custom dataset there are 45 intents and 9 entities.

   - Intents are classes to which each FAQ belongs to and entities are subjects in them.

   - There are a total of 250+ samples.

   - This file is used to train the NLU module which is responsible for intent classification and entity extraction.

# Parts of the Dataset(cont.)

2. *domain.yml*
   - **This file contains all the corresponding responses of each FAQ.**
   - **The responses are classified into 117 different response types.**
   - **There are over 150+ responses arranged in our custom dataset.**
   - **This file also requires all the intents and entities in the nlu.yml file grouped together.**
   - **This file is used to train the Core module used for dialogue management.**

# Parts of the Dataset(cont.)

3. *stories.yml*
   - This file contains 139 stories each mimicking a conversation.

   - Each story consists of an intent defined in the nlu.yml followed by a response defined in the domain.yml.

   - The intent in the story represents a class of query that may come from a user.

   - Each action defines the response our chatbot should give to that input intent.

   - This file is also used to train the Core module used for dialogue management

```
430    - story: 17.4 website
431      steps:
432        - intent: website
433        - action: utter_website
434        - intent: goodbye
435        - action: utter_goodbye
436
437    - story: 18.3 mail_address
438      steps:
439        - intent: mail_address
440        - action: utter_mail_address
441
442    - story: 18.4 mail_address
443      steps:
444        - intent: greet
445        - action: utter_greet
446        - intent: mail_address
447        - action: utter_mail_address
448
449    - story: 19.3 year_of_starting
450      steps:
451        - intent: year_of_starting
452        - action: utter_year_of_starting
453
454    - story: 19.4 year_of_starting
455      steps:
456        - intent: greet
457        - action: utter_greet
458        - intent: year_of_starting
459        - action: utter_year_of_starting
460
```

# Machine Learning Modeling

# Rasa is made up of two separate decoupled units

- NLU module: is used to classify the <u>intent</u> of a given sentence and extract the <u>entities</u>
- Core module: is responsible for the dialogue management of the system

# So how does the entire model work?

- From the input sentence the intents and entities are identified
- The intent and entities along with the previous state are used to update the current state of the conversation
- Policies use the output of the tracker to select an appropriate response from the domain file

**Fig no: How an output message is generated from an input message**

# The NLU module

- It is formed of a pipeline consisting of a series of steps executed consecutively

- The input text is tokenized

- Sparse featurizers like Lexical Syntactic Featurizer, Count Vector Featurizer extract sparse features

- Using these features the intent and entity predictions are made using the DIETClassifer

- The DIETClassifer takes in the features as input and outputs the intents and entities

```
 1  language: bn
 2
 3  ## Pipline 1 = vanilla pipeline
 4  pipeline:
 5    - name: WhitespaceTokenizer
 6    - name: RegexFeaturizer
 7    - name: LexicalSyntacticFeaturizer
 8    - name: CountVectorsFeaturizer
 9    - name: CountVectorsFeaturizer
10      analyzer: char_wb
11      min_ngram: 1
12      max_ngram: 4
13    - name: DIETClassifier
14      epochs: 200
15
```

**Fig no: How an NLU model works**

# How can we customize the NLU module?

- We can add language specific tokenizers for Bangla.

- We can leverage pre-trained dense features on Bangla corpus.

- Concatenating pre-trained dense features along with sparse features highly improves performance.

- FastText and Spacy both provide pre-trained dense features for Bangla.

- BERT provides a language agnostic pre-trained dense feature as well.

# Different customized NLU pipeline



```
94   pipeline:
95   - name: WhitespaceTokenizer
96   - name: CountVectorsFeaturizer
97   - name: CountVectorsFeaturizer
98     analyzer: char_wb
99     min_ngram: 1
100    max_ngram: 4
101  - name: ftfeat.FastTextFeaturizer
102    cache_dir: 'F:/Bot/fastText'
103    file: 'cc.bn.300.bin'
104  - name: DIETClassifier
105    epochs: 100
106
```

```
4   pipeline:
5   - name: custom_tokenizer.BanglaTokenizer
6   - name: RegexFeaturizer
7   - name: LexicalSyntacticFeaturizer
8   - name: CountVectorsFeaturizer
9   - name: CountVectorsFeaturizer
10    analyzer: char_wb
11    min_ngram: 1
12    max_ngram: 4
13  - name: DIETClassifier
14    epochs: 200
15
```
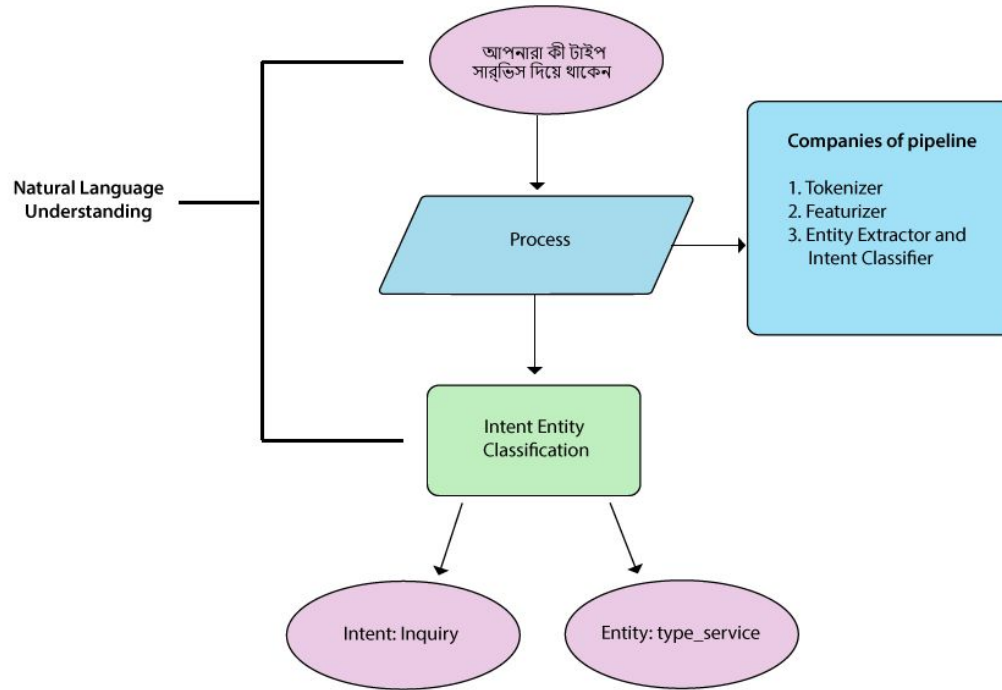
Fig no: pipeline with dense bangla fastText features

Fig no: pipeline with custom bangla tokenizer

# Different customized NLU pipeline(contd.)

```
72    pipeline:
73  - - name: SpacyNLP
74      model: "bn_core_news_sm"
75    - name: SpacyTokenizer
76    - name: SpacyEntityExtractor
77  - - name: SpacyFeaturizer
78      pooling: mean
79  - - name: CountVectorsFeaturizer
80      analyzer: char_wb
81      min_ngram: 1
82      max_ngram: 4
83  - - name: DIETClassifier
84      epochs: 100
85  - - name: FallbackClassifier
86      threshold: 0.2
87      ambiguity_threshold: 0.1
88
```

```
54  - pipeline:
55    - name: "LanguageModelTokenizer"
56  - - name: LanguageModelFeaturizer
57      model_name: "bert"
58      model_weights: "rasa/LaBSE"
59      cache_dir: null
60    - name: RegexFeaturizer
61    - name: LexicalSyntacticFeaturizer
62  - - name: CountVectorsFeaturizer
63      analyzer: char_wb
64      min_ngram: 1
65      max_ngram: 4
66  - - name: DIETClassifier
67      epochs: 500
68  - - name: FallbackClassifier
69      threshold: 0.3
70      ambiguity_threshold: 0.1
71
```

Fig no: pipeline with dense bangla spacy features

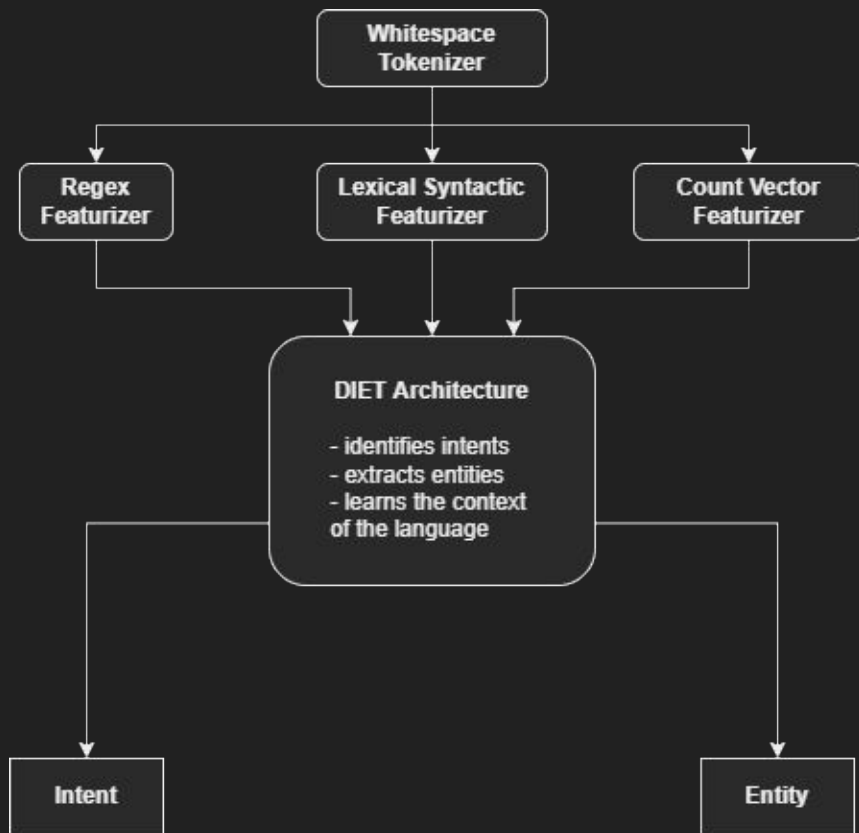Fig no: pipeline with dense language agnostic Bert features

# The Core module

- It includes a tracker and some Policies.

- The tracker keeps track of the conversation state and updates it.

- Policies are of two types: 1) ML based 2) Rule based

- The policies we used for response selection are
    - TED Policy: Transformer based policy to predict the next action to an input text

    - Memoization Policy: A rule-based policy that checks if the current conversation matches our defined stories and responds accordingly

    - Rule Policy: Used to implement fallback response when a out of scope input is fed to the model

# Ablation Study

# Pipeline P1

→ Add

→ Replace

→ Delete

| NLU Pipeline | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| P1 | 75.47 | 63.65 | 75.47 | 67.75 |

# Pipeline P2



| NLU Pipeline | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| P1 | 75.47 | 63.65 | 75.47 | 67.75 |
| P2 | 77.36 | 68.24 | 77.36 | 70.82 |

# Pipeline P3

→ Add

→ Replace

→ Delete

| NLU Pipeline | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| P1 | 75.47 | 63.65 | 75.47 | 67.75 |
| P2 | 77.36 | 68.24 | 77.36 | 70.82 |
| P3 | 77.36 | 66.45 | 77.36 | 70.48 |

# Pipeline P4

Add
Replace
Delete

| NLU Pipeline | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| P1 | 75.47 | 63.65 | 75.47 | 67.75 |
| P2 | 77.36 | 68.24 | 77.36 | 70.82 |
| P3 | 77.36 | 66.45 | 77.36 | 70.48 |
| P4 | 73.58 | 62.74 | 73.58 | 66.49 |

# Pipeline P5

Add

Replace

Delete

| NLU Pipeline | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| P1 | 75.47 | 63.65 | 75.47 | 67.75 |
| P2 | 77.36 | 68.24 | 77.36 | 70.82 |
| P3 | 77.36 | 66.45 | 77.36 | 70.48 |
| P4 | 73.58 | 62.74 | 73.58 | 66.49 |
| P5 | 79.25 | 71.38 | 79.25 | 73.46 |

# Pipeline P6



Custom Bangla Tokenizer

fastText Featurizer | Lexical Syntactic Featurizer | Count Vector Featurizer

DIET Architecture

- identifies intents
- extracts entities
- learns the context of the language

Fallback Classifier

Entity Synonym Mapper

Intent

Entity

→ Add
→ Replace
→ Delete

| NLU Pipeline | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| P1 | 75.47 | 63.65 | 75.47 | 67.75 |
| P2 | 77.36 | 68.24 | 77.36 | 70.82 |
| P3 | 77.36 | 66.45 | 77.36 | 70.48 |
| P4 | 73.58 | 62.74 | 73.58 | 66.49 |
| P5 | 79.25 | 71.38 | 79.25 | 73.46 |
| P6 | 81.13 | 73.27 | 81.13 | 75.32 |

# Pipeline P7



Add

Replace

Delete

| NLU Pipeline | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| P1 | 75.47 | 63.65 | 75.47 | 67.75 |
| P2 | 77.36 | 68.24 | 77.36 | 70.82 |
| P3 | 77.36 | 66.45 | 77.36 | 70.48 |
| P4 | 73.58 | 62.74 | 73.58 | 66.49 |
| P5 | 79.25 | 71.38 | 79.25 | 73.46 |
| P6 | 81.13 | 73.27 | 81.13 | 75.32 |
| P7 | 79.25 | 75.16 | 79.25 | 75.47 |

# Pipeline P8



Add
Replace
Delete

| NLU Pipeline | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| P1 | 75.47 | 63.65 | 75.47 | 67.75 |
| P2 | 77.36 | 68.24 | 77.36 | 70.82 |
| P3 | 77.36 | 66.45 | 77.36 | 70.48 |
| P4 | 73.58 | 62.74 | 73.58 | 66.49 |
| P5 | 79.25 | 71.38 | 79.25 | 73.46 |
| P6 | 81.13 | 73.27 | 81.13 | 75.32 |
| P7 | 79.25 | 75.16 | 79.25 | 75.47 |
| P8 | **83.02** | **80.82** | **83.02** | **80** |

# Experimental Setup

- We split the data into 80-20 train/test split

- THE NLU MODULE
  - Learning rate: 0.5
  - Optimizer: Adam
  - Epochs: 500
  - Minimum ngram for Count Vector Featurizer(sparse featurizer): 1
  - Maximum ngram for Count Vector Featurizer(sparse featurizer): 4
  - Fallback Classifier threshold: 0.3

- THE CORE MODULE
  - Max history: 5
  - Epochs: 200

# Conclusion and Future Work

- Extend and improve the quality of our dataset.

- Use state of the art custom components like language specific dense featurizers for example: BERT embeddings trained on Bangla corpus.

- Change the model in a way as to handle conversations in several languages.

# End-to-End Interaction

# References