



# **TRADITIONAL METHODS FOR HINDI NAMED ENTITY RECOGNITION**

## **TEAM MEMBERS:**

**Akshara Choutapally(2020IMT-024)**

**Deepjyoti Boro(2020IMT-025)**

**Devanshu Patidar(2020IMT-026)**

**Presented to Dr. Sunil Kumar**

# TABLE OF CONTENT

1. INTRODUCTION
2. DATASET
3. METHODOLOGY
4. PERFORMANCE EVALUATION
5. CONCLUSION
6. REFERENCES



# Named Entity Recognition

- NER is a vital NLP task used to **extract and categorize named entities into predefined classes** such as person, location, organization, numeral, and temporal entities.[1]
- **Some Applications of NER:**
  - Information Retrieval
  - Question Answering
  - Sentimental Analysis
- Enables better understanding of text context and semantics.

Cristiano Ronaldo dos Santos Aveiro GOIH ComM  
(Portuguese pronunciation: [kɾiʃ'tjɐnu ãõ'naɫdu]; born 5 February 1985) is a Portuguese professional footballer who plays as a

Model ? English - en\_core\_web\_sm (v3.1.0)

Entity labels (select all)

PERSON  NORP  
 ORG  GPE  LOC  
 PRODUCT  EVENT  
 WORK OF ART  LANGUAGE  
 DATE  TIME  
 PERCENT  MONEY  
 QUANTITY  ORDINAL  
 CARDINAL

Cristiano Ronaldo PERSON dos Santos Aveiro PERSON GOIH ComM (Portuguese pronunciation: [kɾiʃ'tjɐnu ãõ'naɫdu PERSON]; born 5 February 1985 DATE) is a Portuguese professional footballer who plays as a forward for Premier League club Manchester United and captains the Portugal national team. He has won 32 trophies in his career, including seven league titles, five UEFA Champions Leagues, one UEFA European Championship EVENT, and one UEFA Nations League. Ronaldo PERSON holds the records for most appearances (183), most goals (140), and assists (42) in the Champions League, most goals in the European Championship EVENT (14), most international goals by a male player (115), and most international appearances by a European male (186).

# SOME GLIMPSE OF HINER DATASET

datasets

Python

```
DatasetDict({  
    train: Dataset({  
        features: ['id', 'tokens', 'ner_tags'],  
        num_rows: 75827  
    })  
    validation: Dataset({  
        features: ['id', 'tokens', 'ner_tags'],  
        num_rows: 10851  
    })  
    test: Dataset({  
        features: ['id', 'tokens', 'ner_tags'],  
        num_rows: 21657  
    })  
})
```

# SOME GLIMPSE OF HINER DATASET

```
datasets["train"][2]

{'id': '2',
'tokens': ['रामनगर',
'इगलास',
',',
'अलीगढ़',
',',
'उत्तर',
'प्रदेश',
'स्थित',
'एक',
'गाँव',
'है'],
'ner_tags': [0, 0, 6, 0, 6, 0, 3, 6, 6, 6, 6]}
```

Python

# SOME GLIMPSE OF HINER DATASET

```
label_list = datasets["train"].features[f"ner_tags"].feature.names  
label_list
```

Python

```
[ 'O', 'B-PER', 'I-PER', 'B-LOC', 'I-LOC', 'B-ORG', 'I-ORG' ]
```

'O': Outside of any named entity

'B-PER': Beginning of a person's name

'I-PER': Inside of a person's name (if it consists of multiple tokens)

'B-LOC': Beginning of a location name

'I-LOC': Inside of a location name (if it consists of multiple tokens)

'B-ORG': Beginning of an organization name

'I-ORG': Inside of an organization name (if it consists of multiple tokens)

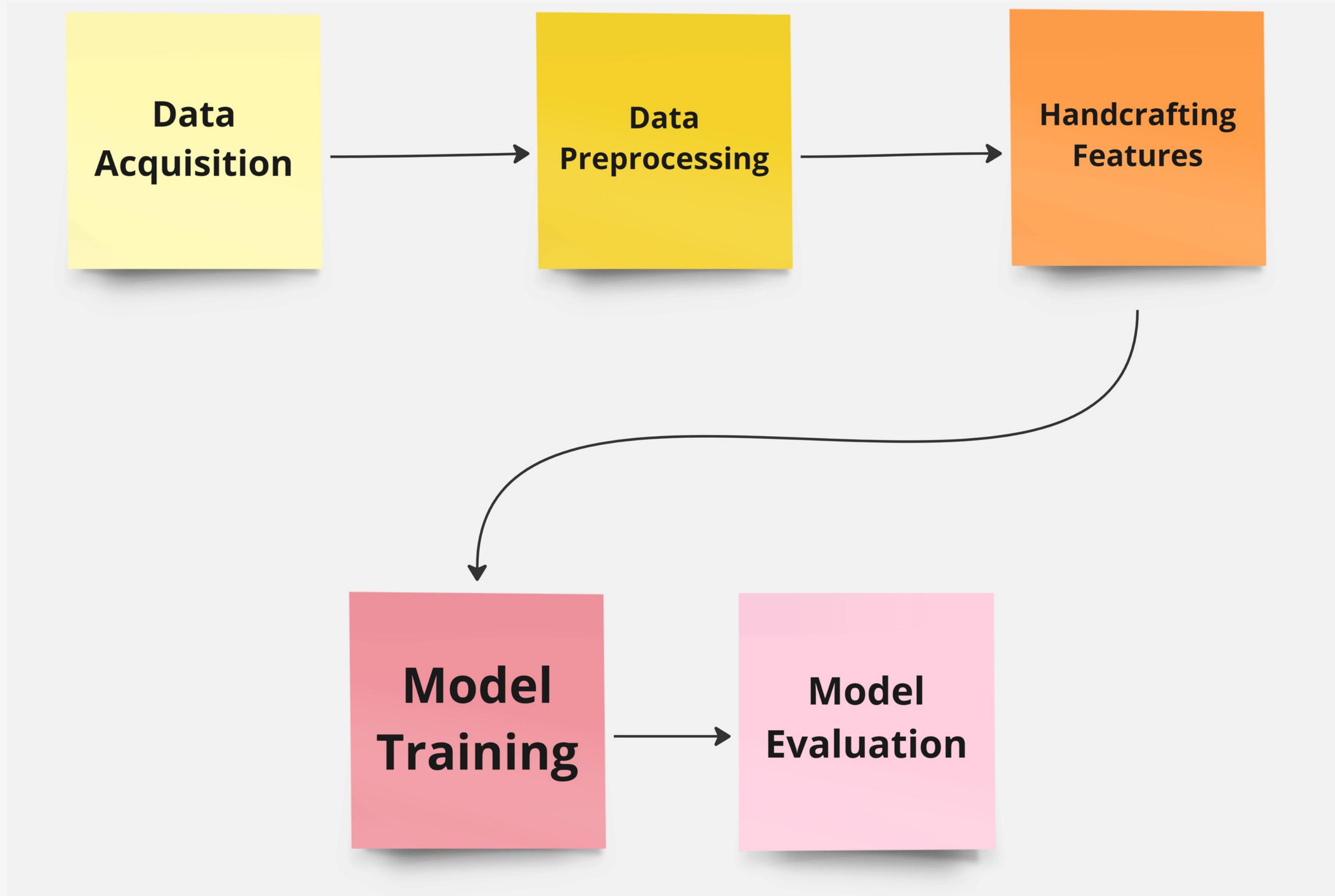
# SOME GLIMPSE OF HINER DATASET

```
show_random_elements(datasets["train"])
```

Python

	<b>id</b>	<b>tokens</b>	<b>ner_tags</b>
0	61733	[जल्दी, ही, रामनरेश, सरवन, भी, 5, रन, बना, कर, पैवेलियन, लौट, गए, जब, सरवन, के, रूप, में, पाँचवां, विकेट, गिरा, वेस्टइंडीज़, के, 12वें, ओवर, में, मात्र, 64, रन, बने, थे, .]	[I-ORG, I-ORG, I-PER, B-ORG, I-ORG, O, I-ORG, I-ORG]
1	59780	[मेरा, साथी, 1985, में, बनी, हिन्दी, भाषा, की, फ़िल्म, है।]	[I-ORG, I-ORG, I-ORG]
2	59142	[१४४, किलोमीटर, लंबा, यह, राजमार्ग, थंजावुर, को, मनमदुरई, से, जोड़ता, है।]	[I-ORG, I-ORG, I-ORG, I-ORG, I-ORG, O, I-ORG, O, I-ORG, I-ORG, I-ORG]
3	65610	[जब, बनार्ड, मैडॉफ़, ने, सभी, 11, आरोप, स्वीकार, कर, लिए, तो, न्यायाधीश, ने, इसे, मंजूरी, दे, दी, .]	[I-ORG, I-PER, B-ORG, I-ORG, I-ORG]
4	75660	[ज्योया, अफ़रोज़, (जन्म: १०, जनवरी, १९९४, ), एक, भारतीय, अभिनेत्री, और, मॉडल, हैं।]	[I-PER, B-ORG, I-ORG, I-ORG]
5	17302	[इंतज़ाम, .]	[I-ORG, I-ORG]

# METHODOLOGY



# DATA ACQUISITION

- **HiNER dataset** available from the **`datasets` library** was used.[2]
- This dataset is specifically designed for NER tasks in Hindi and includes a variety of annotated named entities.
- The dataset was divided into a training set and a test set, with the training set used to train the models and the test set used to evaluate their performance.

# DATA PREPROCESSING

- Data preprocessing is crucial to prepare the raw data for effective machine learning. The preprocessing steps included:
- **Tokenization:** Splitting the text data into individual words or tokens. Tokenization is essential for creating features from the text data.
- **Part-of-Speech (POS) Tagging:** Each token in the dataset was annotated with a part-of-speech tag **using the StanfordNLP library** configured for Hindi. POS tags provide syntactic information about each word, which is valuable for understanding its role in a sentence.  
[3]
- **Data Limiting:** Due to computational constraints or to balance the dataset, the number of sentences used from the dataset was limited to the **first 2000 entries for both training and testing phases**.

# FEATURE ENGINEERING

- Feature engineering is a critical step that involves creating input features for the machine learning models from the raw data. The features extracted for this project included:
- **POS Tags:** The POS tag of each word was used as a feature because different types of words (nouns, verbs, adjectives) tend to have different likelihoods of being named entities.
- **Binary Indicators:** Features indicating whether a word is at the beginning ('BOS') or end ('EOS') of a sentence were included. These features help the model recognize boundary-related patterns that are common in named entities.
- **Contextual Features:** The **words immediately before and after a given word were used as features.** It's important in NER tasks as the likelihood of a word being a named entity can depend heavily on its surrounding words.

# MODEL TRAINING AND EVALUATION

- With the features prepared, the next step was to train different machine learning models. Three models were chosen for this study:[4]

**Decision Tree Classifier**

**Random Forest Classifier**

**Support Vector Machine(SVM)**

- Each model was trained on the extracted features and were tested using the unseen test set.
- The **evaluation metrics** used were **precision, recall, and F1-score**
- Confusion matrices** were generated for each model to visually assess how well each model was performing across different classes of named entities.

# **PERFORMANCE EVALUATION**

# PERFORMANCE METRICS USED

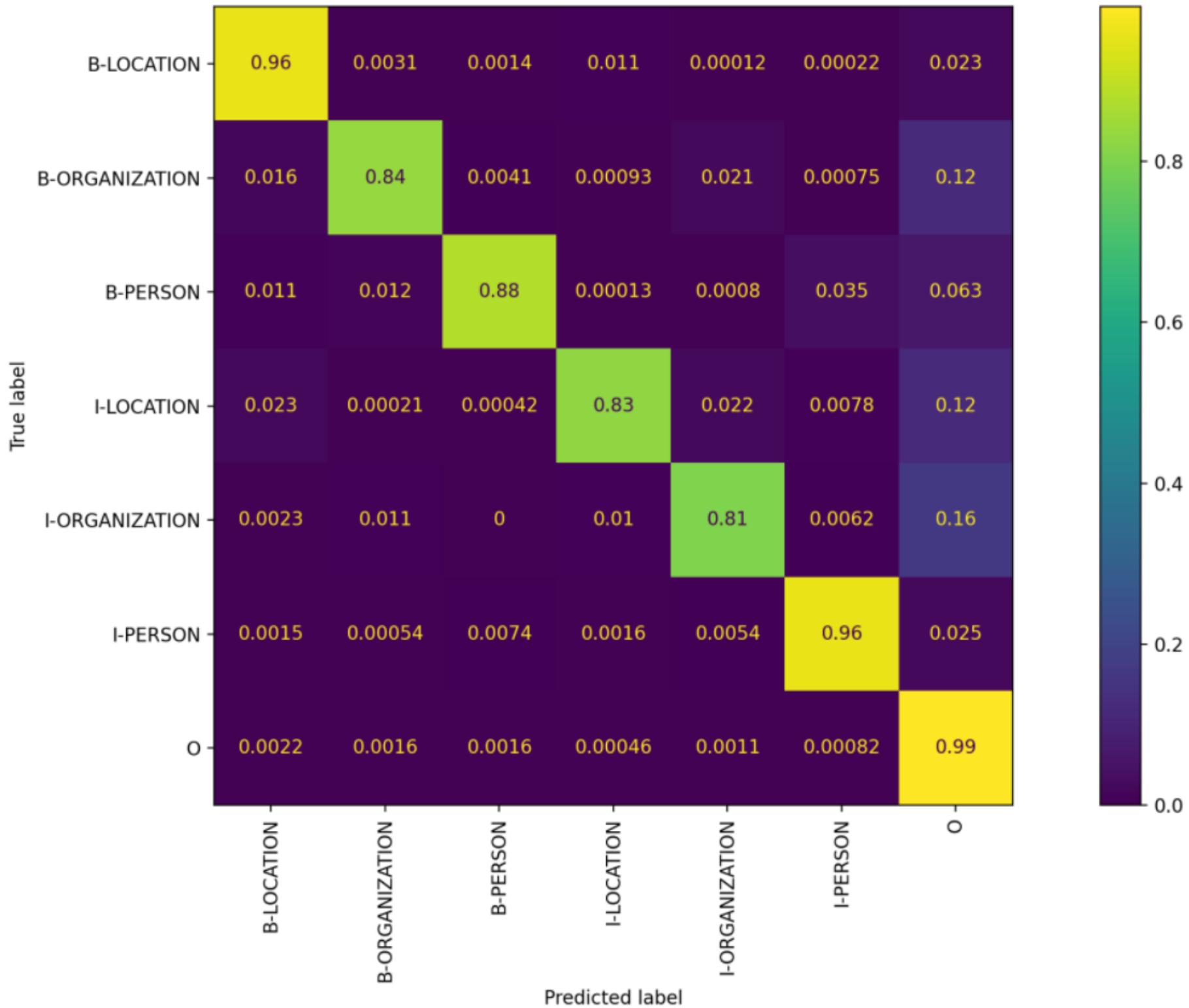
$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

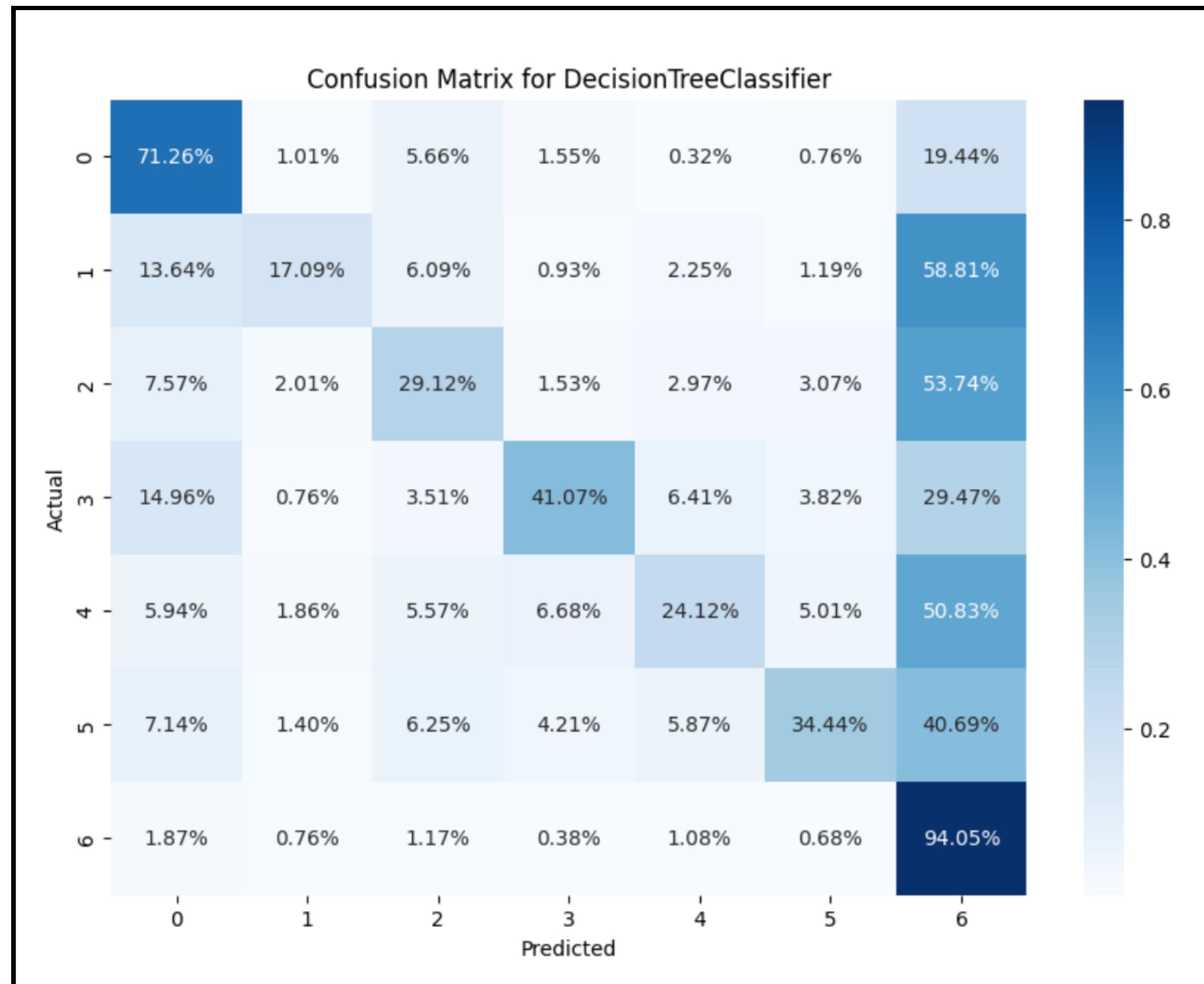
$$F1 = \frac{2 \times Precision \times Recall}{Precision + Recall}$$

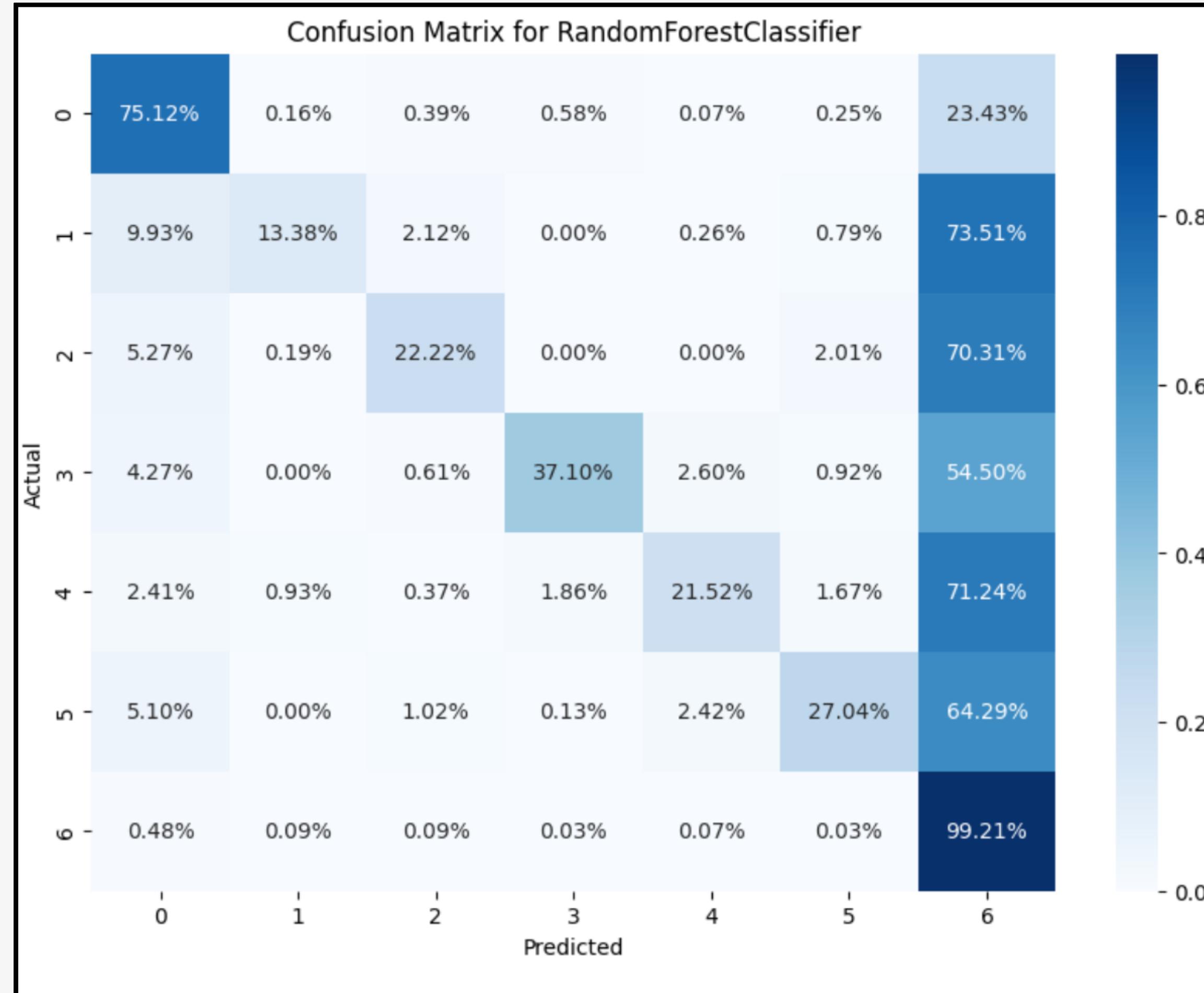
MODEL	Weighted F1 score	Macro Avg F1 score
<b>Base Paper Model</b>	<b>0.92</b>	<b>0.86</b>
Random Forest	0.91	0.49
Decision Tree	0.89	0.43
SVM	0.89	0.35

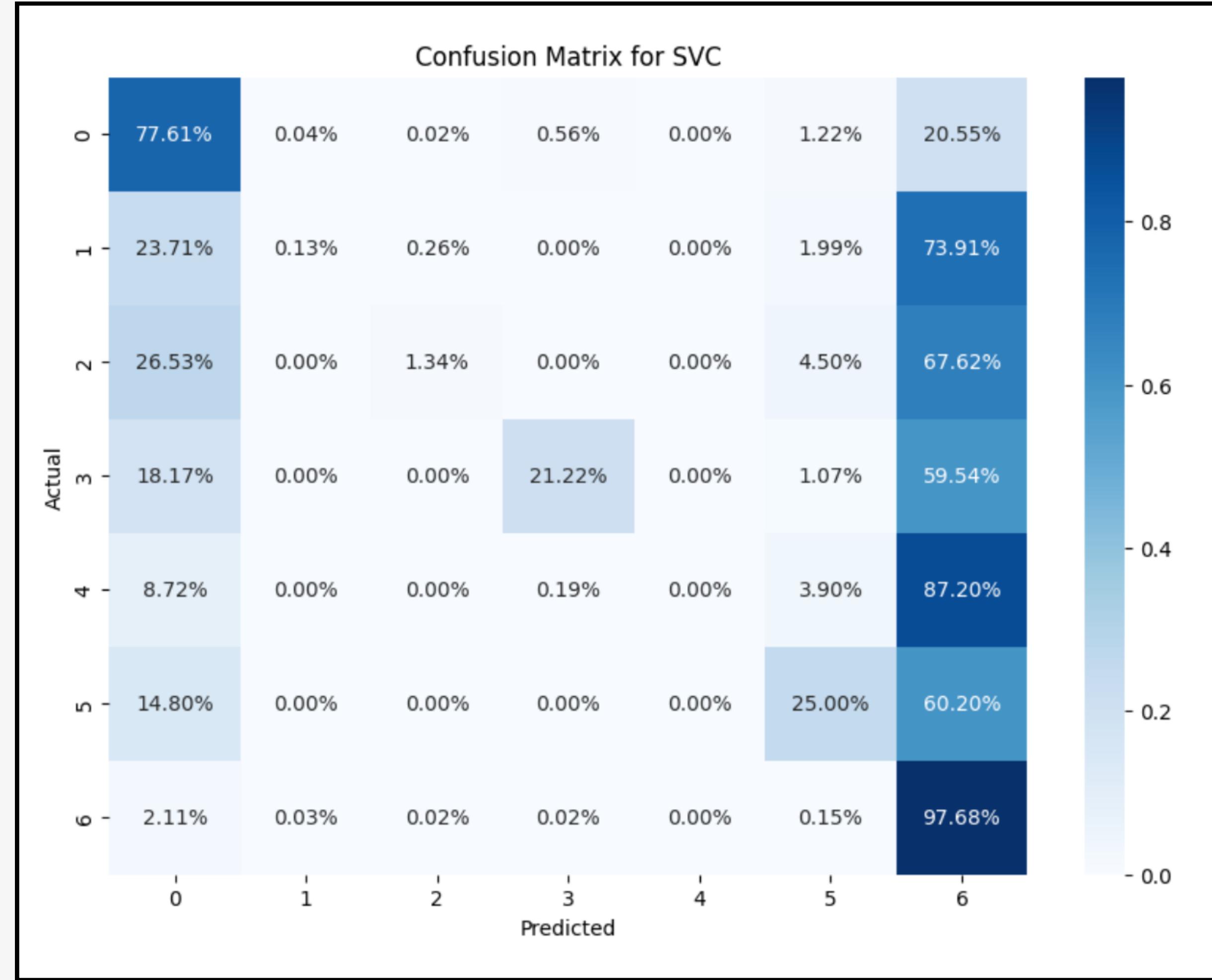
# BASE PAPER RESULTS



# OUR RESULTS:







Classification Report for DecisionTreeClassifier:

		precision	recall	f1-score	support
	0	0.81	0.71	0.76	3778
	1	0.20	0.17	0.19	508
	2	0.26	0.24	0.25	671
	3	0.35	0.38	0.36	431
	4	0.26	0.27	0.26	349
	5	0.22	0.31	0.26	495
	6	0.95	0.95	0.95	34477
	accuracy			0.89	40709
	macro avg	0.44	0.43	0.43	40709
	weighted avg	0.89	0.89	0.89	40709

Classification Report for RandomForestClassifier:

	precision	recall	f1-score	support
0	0.90	0.74	0.81	3778
1	0.65	0.14	0.23	508
2	0.58	0.20	0.30	671
3	0.82	0.35	0.49	431
4	0.64	0.19	0.29	349
5	0.62	0.23	0.33	495
6	0.93	0.99	0.96	34477
accuracy			0.92	40709
macro avg	0.73	0.41	0.49	40709
weighted avg	0.91	0.92	0.91	40709

## Classification Report for SVC:

	precision	recall	f1-score	support
0	0.71	0.78	0.74	3778
1	0.08	0.00	0.00	508
2	0.50	0.02	0.05	671
3	0.92	0.23	0.37	431
4	0.00	0.00	0.00	349
5	0.43	0.25	0.32	495
6	0.93	0.98	0.95	34477
accuracy			0.91	40709
macro avg	0.51	0.32	0.35	40709
weighted avg	0.88	0.91	0.89	40709

# CONCLUSION

# BASE PAPER MODEL OUTPERFORMED TRADITIONAL METHODS BY LARGE MARGIN.

why?

- Limited training data as a consequence of limited computational resources.
- Traditional methods not suitable to handle sequential data and patterns.
- Traditional methods rely on hand-crafted features, while modern models can automatically learn representations from raw text data.
- Lack of tools and library support for regional languages.

## FUTURE SCOPE

- Use of models that capture temporal dependencies in the text data.
- Use of Models that understand Contextual nature of the tokens.
- Improving feature engineering , adding morphological features.
- Generating text classification datasets, sentiment analysis datasets from this dataset.

# REFERENCES

---

- [1] A. Awan, “What is Named Entity Recognition (NER)?,” Datacamp.  
<https://www.datacamp.com/blog/what-is-named-entity-recognition-ner>
- [2] R. Murthy, P. Bhattacharjee, R. Sharnagat, J. Khatri, D. Kanojia, and P. Bhattacharyya, “HiNER: A Large Hindi Named Entity Recognition Dataset.” Accessed: Apr. 09, 2024. [Online]. available:<https://www.cse.iitb.ac.in/~pb/papers/lrec22-ner.pdf>
- [3] “stanfordnlp/stanza,” GitHub, <https://github.com/stanfordnlp/stanza/>
- [4] P. Bose, S. Srinivasan, W. C. Sleeman, J. Palta, R. Kapoor, and P. Ghosh, “A Survey on Recent Named Entity Recognition and Relationship Extraction Techniques on Clinical Texts,” *Applied Sciences*, vol. 11, no. 18, p. 8319, Sep. 2021, doi: <https://doi.org/10.3390/app11188319>.

# THANK YOU

