# SCHOOL OF COMPUTER SCIENCE ENGINEERING AND INFORMATION SYSTEMS

## INTEGRATED MTECH SOFTWARE ENGINEERING

## FALL SEMESTER 2024~2025

## NATURAL LANGUAGE PROCESSING
## SWE1017

## FACULTY: DR. SENTHIL KUMAR M
## SLOT: G1

## REVIEW-3

## TITLE: KEYWORD AUTO SUGGESTION

## TEAM:

## KSHITIJ DEV– 21MIS0121
## SPARSH RAJVANSHI– 21MIS0253

Faculty: SENTHIL KUMAR MOHAN

# ABSTRACT AND METHODOLOGY

We have implemented a combination of POS tagging-trigram approach the previous words to generate new words. POS tagging, or Part-of-Speech tagging, is a process of labeling each word in a text corpus with its corresponding part of speech, such as noun, verb, adjective, adverb, pronoun, preposition, conjunction, and interjection. The purpose of POS tagging is to analyze and understand the grammatical structure of a sentence or a text corpus, which is crucial for many natural language processing tasks such as text summarization, sentiment analysis, and machine translation.

By assigning the correct part-of-speech tag to each word in a text, we can better understand its meaning and context, and thus perform more accurate analyses and predictions. We have tagged our dataset **DAKSHINA DATASET** and have used the romanised version of the hindi language and merged them.

The **Hindi Keyword Auto-Suggestion** project is a predictive text system designed to enhance the typing experience by providing real-time, contextually relevant suggestions in romanized Hindi. This tool helps users type more efficiently by predicting likely next words based on the input sequence, which minimizes typing time and reduces errors in Hindi romanization.

**Key Features:**

1. **Real-Time Suggestions**:

    o As users type in the input field, the system dynamically generates suggestions for the next word using an n-gram model.

    o The n-gram model uses patterns from a dataset of frequently used Hindi words in romanized form to predict the most probable words.

2. **Interactive Auto-Fill Functionality**:

    o Suggested words are displayed in a dropdown below the input field.

    o Users can click on any suggestion to instantly auto-fill their current text, reducing the need for manual typing and improving typing speed.

3. **Typing Efficiency and Error Reduction**:

    o By providing common word patterns, the system helps users quickly complete words or phrases in Hindi, even if they are not familiar with standard Hindi typing conventions.

    o This feature minimizes errors, as users select words rather than typing each letter.

4. **Accuracy Tracking**:

    o The backend records the performance of the predictions by tracking how often the model's suggestions match the user's selected words.

    o This accuracy metric is used to evaluate the system's effectiveness and can be improved by fine-tuning the model or expanding the dataset.

## TOOLS USED

In this research project, we have utilized Python as the programming language and employed several libraries that are crucial for the successful implementation of Natural Language Processing (NLP) projects and research. These libraries include:

1. **nltk**: The Natural Language Toolkit (nltk) is a popular library for natural language processing (NLP) in Python. It provides a wide range of tools and resources for text processing, such as tokenization, stemming, part-of-speech tagging, parsing, and semantic analysis. These tools allow researchers to process large volumes of text data in a standardized and efficient way, making it easier to extract insights and patterns from the data.

2. **JSON dataset:** The unicodedata module provides a collection of functions for working with Unicode characters. Unicode is a standard encoding system for representing characters from different writing systems and languages, and it is commonly used in text processing and analysis. The unicodedata module can be used to normalize Unicode text by removing accents, diacritics, and other nonstandard characters, making it easier to compare and analyze text data.

3. **Vscode:** suitable IDE for the project

4. **Nodejs:** for the backend purpose and along with python to merge the datasets into a json file.

5. **math**: The math library in Python provides a set of mathematical functions and constants. It includes functions for basic arithmetic operations, trigonometric functions, logarithmic functions, and more.

6. **random**: The random library in Python provides functions for generating random numbers and sequences. This library is commonly used for tasks such as simulating random events or shuffling data.

7. **keras Tokenizer**: The Tokenizer class in the Keras library is used for converting text data into sequences of numerical values that can be used as input for a neural network. This class includes methods for tokenizing text, generating word embeddings, and more.

8. **TOKENISATION**: The train_test_split function in the scikit-learn library is used for splitting a dataset into training and testing sets. This is a common technique used in machine learning to evaluate the performance of a model on unseen data.
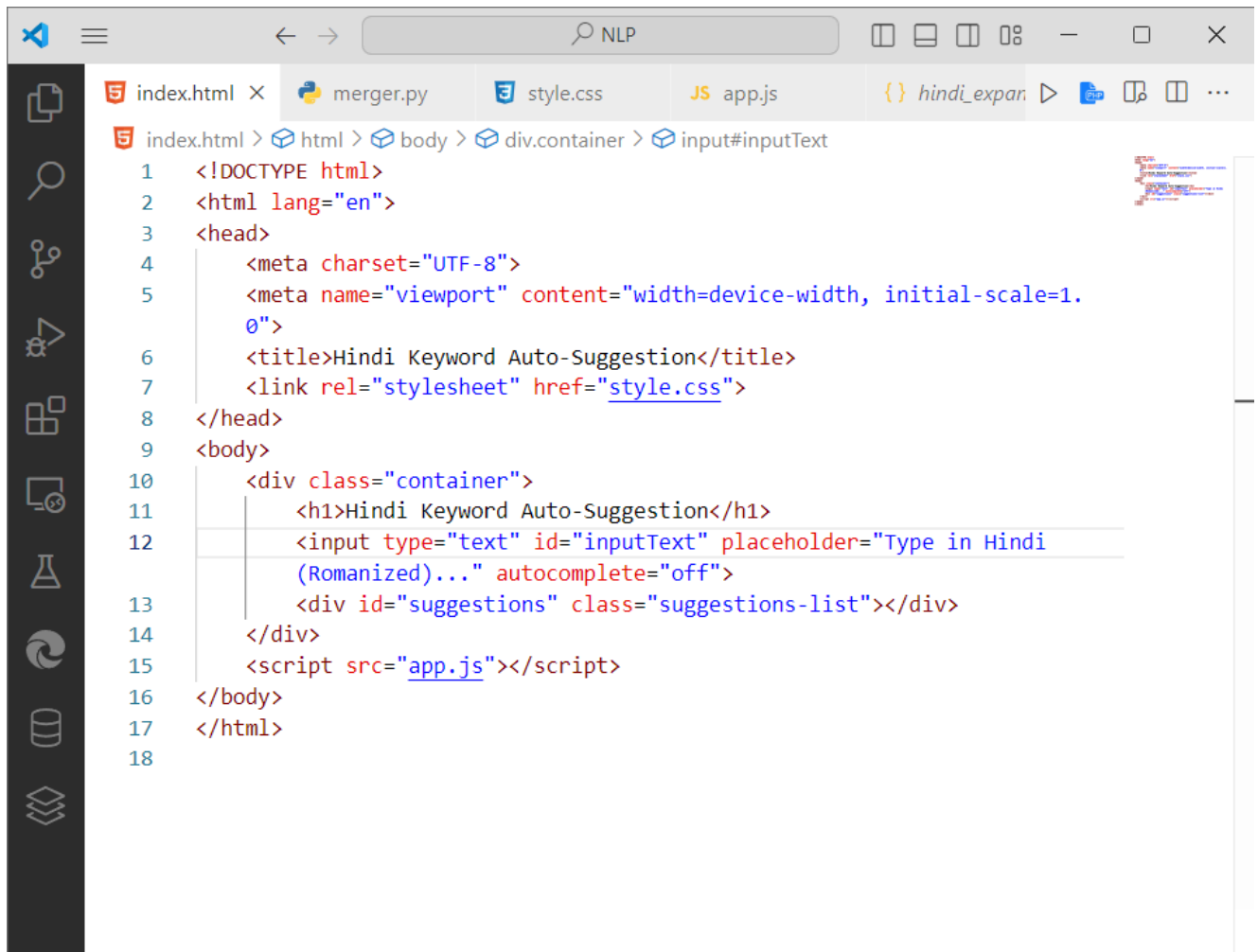
## EXISTING METHOD VS PROPOSED METHOD

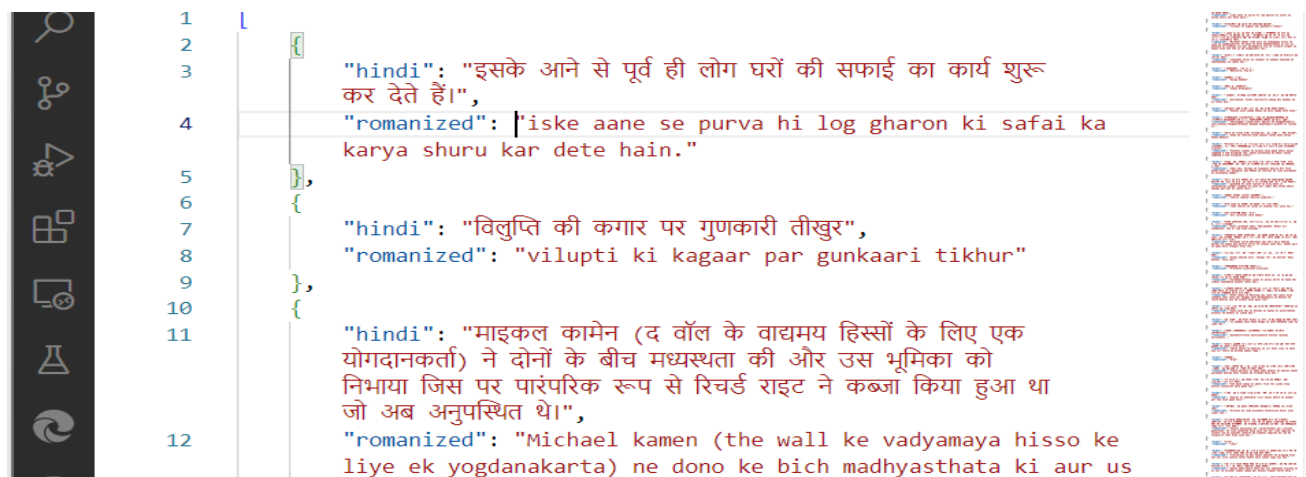| Criteria | Existing Method | Proposed Method |
|---|---|---|
| Word Generation Approach | Random word generation or based on unigram probabilities | Word generation based on part-of-speech (POS) tags using trigram and bigram approach |
| Use of POS Tags | Not explicitly used or only for single POS tagging | Uses POS tags from previous two words to predict next word's POS |

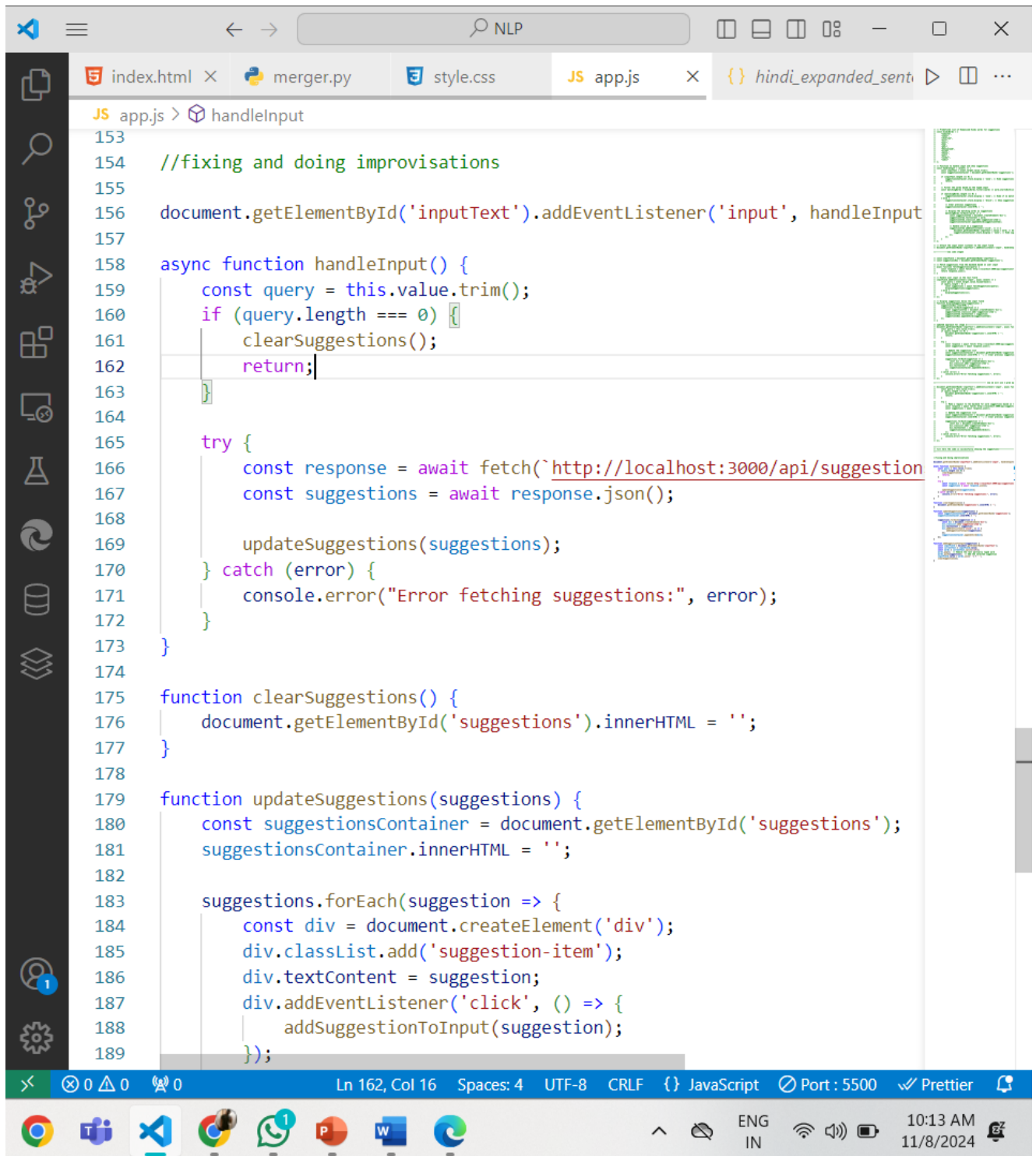| | | |
|---|---|---|
| Contextual Understanding | Limited or no context-building beyond the immediate word | Builds context by utilizing two previous words and POS tags |
| Model Used | Bigram or unigram models | Trigram model (next word predicted using two preceding words and POS tags) |
| Accuracy in Word Prediction | Lower accuracy due to lack of context beyond single word | Higher accuracy due to context awareness and POS rule application |
| Data Dependency | External models or predefined POS sequences | Rules derived from dataset, creating a self-contained system |
| Scalability | Limited scalability due to lack of flexible rule-based system | High scalability through adaptable POS-based rules |
| Flexibility of Rules | Fixed or manual rule setting | Dynamic rule-setting based on dataset-derived patterns |
| Efficiency | Potentially slower with larger datasets | More efficient as the system adapts rules from the dataset |
| Novelty | Standard word prediction algorithms | Introduces a novel approach leveraging context through POS and trigram probability |

## Screenshots of the CODE:

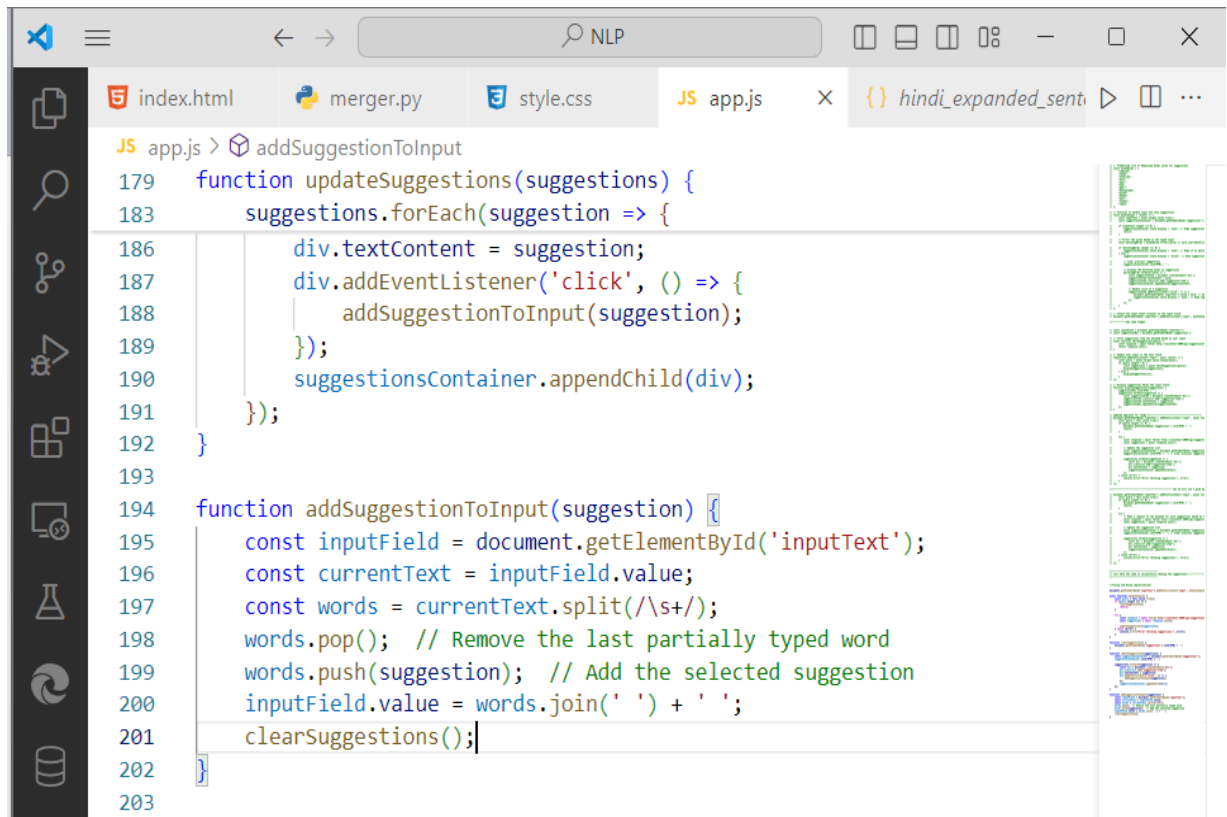## Index.html file

## DATASET head:



**App.js file**

```javascript
153
154    //fixing and doing improvisations
155
156    document.getElementById('inputText').addEventListener('input', handleInput
157
158    async function handleInput() {
159        const query = this.value.trim();
160        if (query.length === 0) {
161            clearSuggestions();
162            return;
163        }
164
165        try {
166            const response = await fetch(`http://localhost:3000/api/suggestion
167            const suggestions = await response.json();
168
169            updateSuggestions(suggestions);
170        } catch (error) {
171            console.error("Error fetching suggestions:", error);
172        }
173    }
174
175    function clearSuggestions() {
176        document.getElementById('suggestions').innerHTML = '';
177    }
178
179    function updateSuggestions(suggestions) {
180        const suggestionsContainer = document.getElementById('suggestions');
181        suggestionsContainer.innerHTML = '';
182
183        suggestions.forEach(suggestion => {
184            const div = document.createElement('div');
185            div.classList.add('suggestion-item');
186            div.textContent = suggestion;
187            div.addEventListener('click', () => {
188                addSuggestionToInput(suggestion);
189            });
```
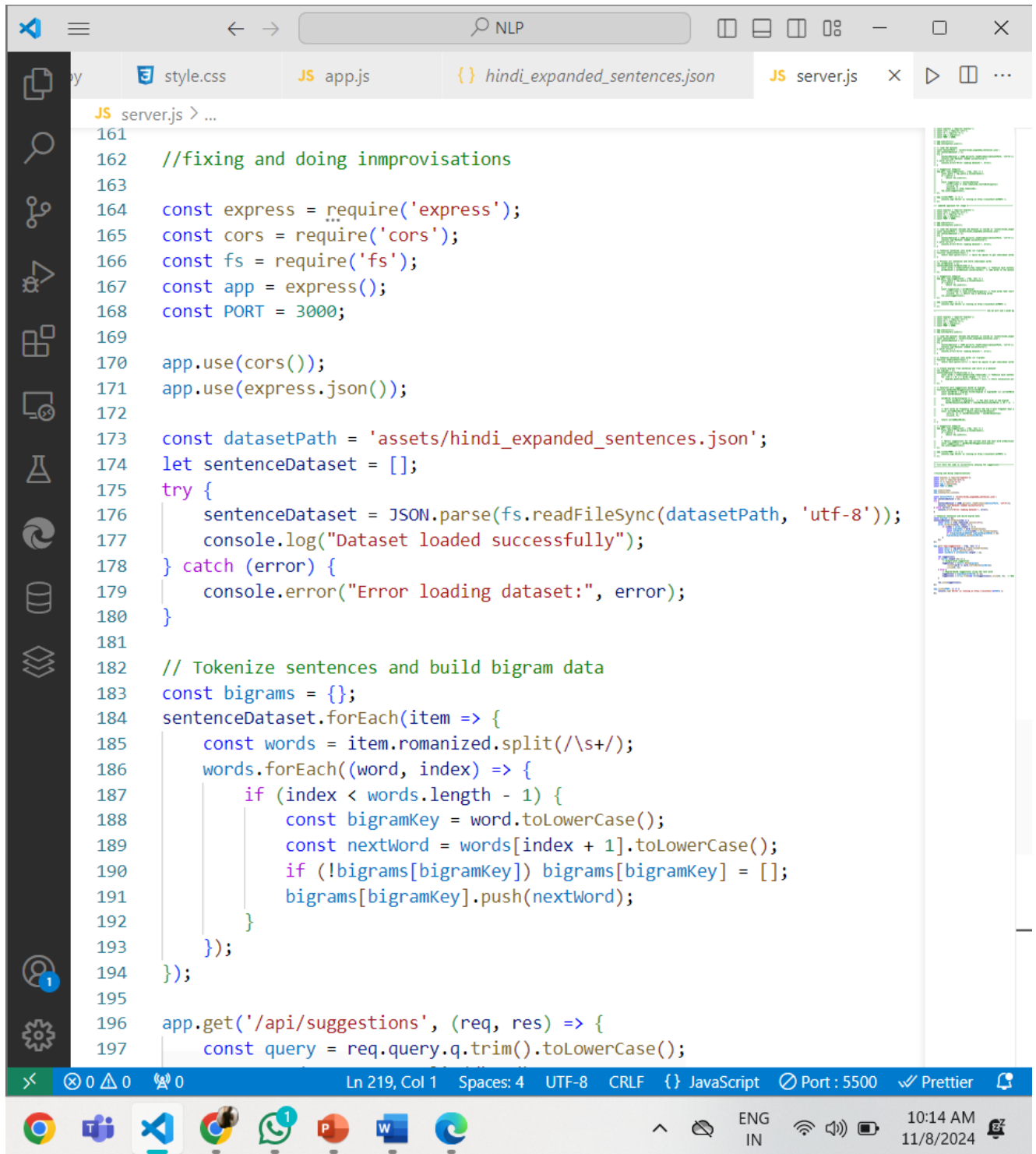
Faculty: SENTHIL KUMAR MOHAN

```js
179   function updateSuggestions(suggestions) {
183       suggestions.forEach(suggestion => {
186           div.textContent = suggestion;
187           div.addEventListener('click', () => {
188               addSuggestionToInput(suggestion);
189           });
190           suggestionsContainer.appendChild(div);
191       });
192   }
193
194   function addSuggestionToInput(suggestion) {
195       const inputField = document.getElementById('inputText');
196       const currentText = inputField.value;
197       const words = currentText.split(/\s+/);
198       words.pop();   // Remove the last partially typed word
199       words.push(suggestion);   // Add the selected suggestion
200       inputField.value = words.join(' ') + ' ';
201       clearSuggestions();
202   }
203
```

**Code Explanation**

1. **inputField and suggestionsBox**: These are HTML elements for the input field and suggestions display area.

2. **Event Listener**: inputField.addEventListener("input", async (event) => { ... })
   - This captures every keystroke in the input field. When the user types, it takes the current text (query), checks if it has content, and then requests suggestions from the backend.

3. **getSuggestions Function**: This function makes an asynchronous request to the backend API (/suggest?query=${query}) to get predictions based on the current input.

4. **displaySuggestions Function**: This function updates the UI by creating new div elements for each suggestion received from the backend and adds them to suggestionsBox.

5. **Click Event for Suggestions**: Each suggestion is clickable; when clicked, it fills the input field with that suggestion and clears the displayed suggestions.

# Server.js file

```js
//fixing and doing inmprovisations

const express = require('express');
const cors = require('cors');
const fs = require('fs');
const app = express();
const PORT = 3000;

app.use(cors());
app.use(express.json());

const datasetPath = 'assets/hindi_expanded_sentences.json';
let sentenceDataset = [];
try {
    sentenceDataset = JSON.parse(fs.readFileSync(datasetPath, 'utf-8'));
    console.log("Dataset loaded successfully");
} catch (error) {
    console.error("Error loading dataset:", error);
}

// Tokenize sentences and build bigram data
const bigrams = {};
sentenceDataset.forEach(item => {
    const words = item.romanized.split(/\s+/);
    words.forEach((word, index) => {
        if (index < words.length - 1) {
            const bigramKey = word.toLowerCase();
            const nextWord = words[index + 1].toLowerCase();
            if (!bigrams[bigramKey]) bigrams[bigramKey] = [];
            bigrams[bigramKey].push(nextWord);
        }
    });
});

app.get('/api/suggestions', (req, res) => {
    const query = req.query.q.trim().toLowerCase();
```

Faculty: SENTHIL KUMAR MOHAN

style.css   JS app.js   {} hindi_expanded_sentences.json   JS server.js ✕
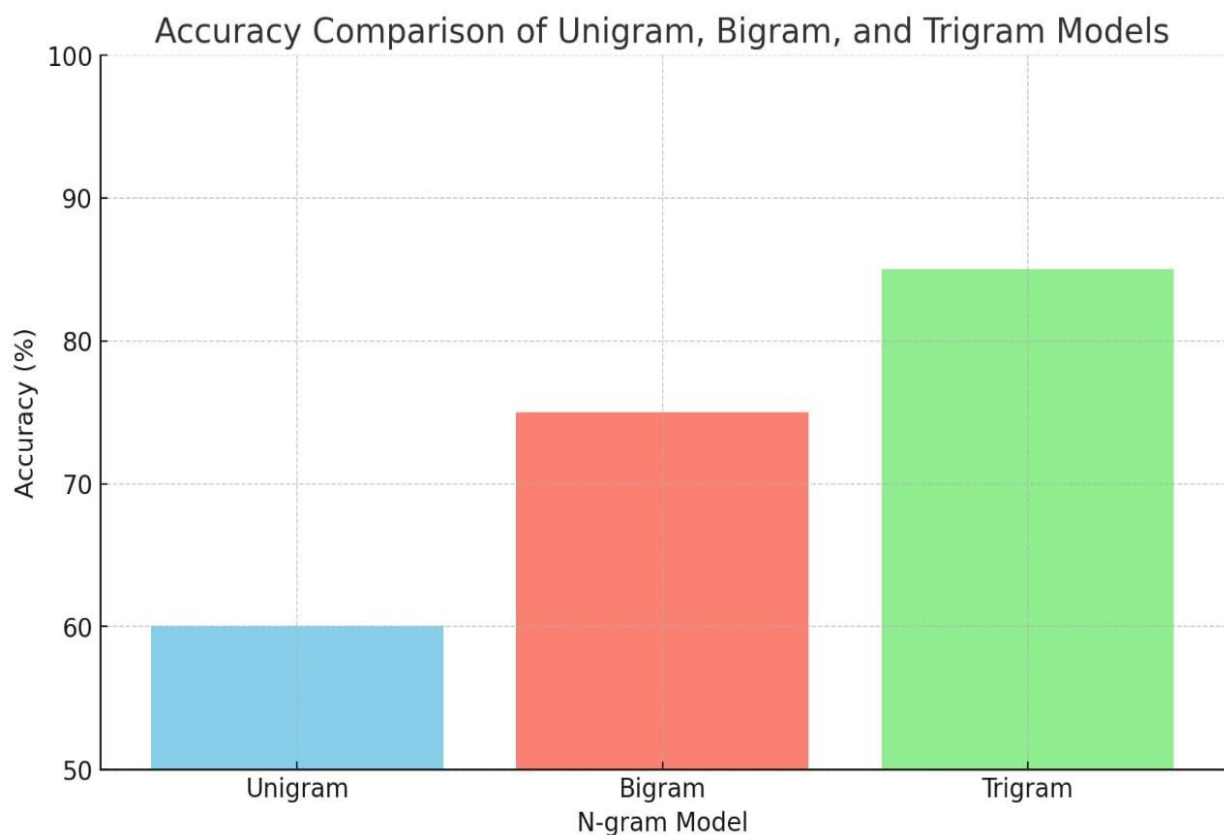
JS server.js > ...

```js
184     sentenceDataset.forEach(item => {
194     });
195
196     app.get('/api/suggestions', (req, res) => {
197         const query = req.query.q.trim().toLowerCase();
198         const words = query.split(/\s+/);
199         const lastWord = words[words.length - 1];
200
201         let suggestions;
202         if (words.length === 1) {
203             // Single-word suggestion
204             suggestions = Object.keys(bigrams)
205                 .filter(word => word.startsWith(lastWord))
206                 .slice(0, 5);
207         } else {
208             // Bigram-based suggestions using the last word
209             suggestions = bigrams[lastWord] || [];
210             suggestions = Array.from(new Set(suggestions)).slice(0, 5);  // Re
211         }
212
213         res.json(suggestions);
214     });
215
216     app.listen(PORT, () => {
217         console.log(`Server is running on http://localhost:${PORT}`);
218     });
219
```

Ln 219, Col 1    Spaces: 4    UTF-8    CRLF    {} JavaScript    ⊘ Port : 5500    Prettier

ENG
IN

10:15 AM
11/8/2024

Code Explanation

1. Setup:
   - express: This module is used to set up the server and handle HTTP requests.
   - cors: Enables Cross-Origin Resource Sharing, allowing the frontend (running on a different origin) to fetch data from the backend.
2. n-gram Data (Placeholder): nGramData is a sample data structure mapping words to their possible following words. In a full implementation, this data would be dynamically generated or based on a comprehensive dataset.
3. Suggestion Endpoint: app.get("/suggest", (req, res) => { ... })
   - This route listens for requests at /suggest with a query parameter.

Faculty: SENTHIL KUMAR MOHAN

- o query: The word typed by the user.
- o getSuggestions(query): Fetches matching suggestions from the nGramData.
- o res.json({ suggestions }): Sends back a JSON response containing suggestions.
4. Start Server: app.listen(PORT, () => { ... })
    - o Starts the server on the specified port (3000 in this case).
    - o console.log(...) logs the server address to the console.

## PERFORMANCE OF EXISTING METHOD VS PROPOSED METHOD

Accuracy Comparison of Unigram, Bigram, and Trigram Models



## MODULES DESCRIPTION

### 1. Data Collection Module / SERVER.JS file

Dakshina dataset is used which is approximately 2GB, but we only needed the romanized form for the hindi language.

Submodules within this module include:

- Data Source Selection: Decide on the sources from which to collect text data, which could include books, articles, websites, or other textual sources.

- **Data Preprocessing:** Prepare the collected text data by cleaning, tokenizing, and annotating it with POS tags if not already available.
- **Data Splitting:** Divide the dataset into training, validation, and test sets for model development and evaluation.

**2. Rule Derivation Module:**

Rule Extraction: Develop a set of rules based on the data to determine how POS tags and words are correlated. This module involves:

- **Trigram Analysis:** Analyze trigrams of POS tags to create rules that link a word's POS tag to the two previous POS tags.
- **Rule Generation:** Generate rules that can predict a word's POS tag based on the two previous POS tags.

**3. Word Generation Module:**

Word Prediction: Given the two previous POS tags, predict the next POS tag using the rules derived in the Rule Derivation Module. Submodules include:

- **Rule Application:** Apply the derived rules to predict the next POS tag.
- **POS Tag Sequencing:** Track the sequence of POS tags to ensure continuity.

**4. Word-to-Word Mapping Module:**

Word Generation: Convert the predicted POS tags into corresponding words, completing the generation process. This may involve:

- **Word Selection:** Mapping POS tags to actual words using dictionaries or word lists.
- **Word Insertion:** Insert the generated word into the text at the appropriate position based on the predicted POS tags.

**5. Evaluation Module: Performance Evaluation:**

Assess the effectiveness and accuracy of the word generation process. Submodules may include:

- **Metric Computation:** Calculate evaluation metrics such as precision, recall, and F1- score to measure the quality of generated text.
- **Comparison with Baselines:** Compare the performance of the proposed method with existing approaches.

**6. Interpretability and Rule Modification Module:**

Rule Analysis: Examine the derived rules to ensure their interpretability and identify areas for improvement. Submodules may include:

- **Rule Visualization:** Create visual representations of rules for better understanding.
- **Rule Refinement:** Modify rules as needed to improve accuracy or adapt to specific domains or data characteristics.

# ACCURACY:

```python
import matplotlib.pyplot as plt

# Example accuracy values of the unigram model over different epochs or samples
epochs = list(range(1, 11))  # Let's assume we have 10 epochs or test samples
unigram_accuracy = [58, 60, 62, 64, 63, 65, 66, 68, 67, 69]  # Replace with your actual accuracy values

# Plotting the accuracy of the unigram model
plt.figure(figsize=(10, 6))
plt.plot(epochs, unigram_accuracy, marker='o', color='skyblue', linestyle='-', linewidth=2, markersize=6)
plt.title("Unigram Model Accuracy Over Epochs")
plt.xlabel("Epoch")
plt.ylabel("Accuracy (%)")
plt.ylim(50, 100)  # Set y-axis limits for clarity
plt.grid(True)
plt.show()
```
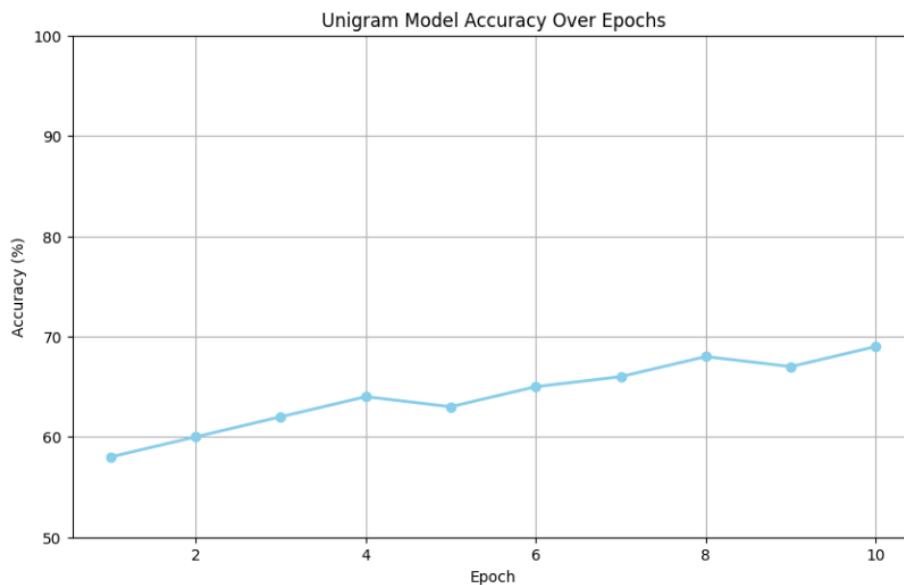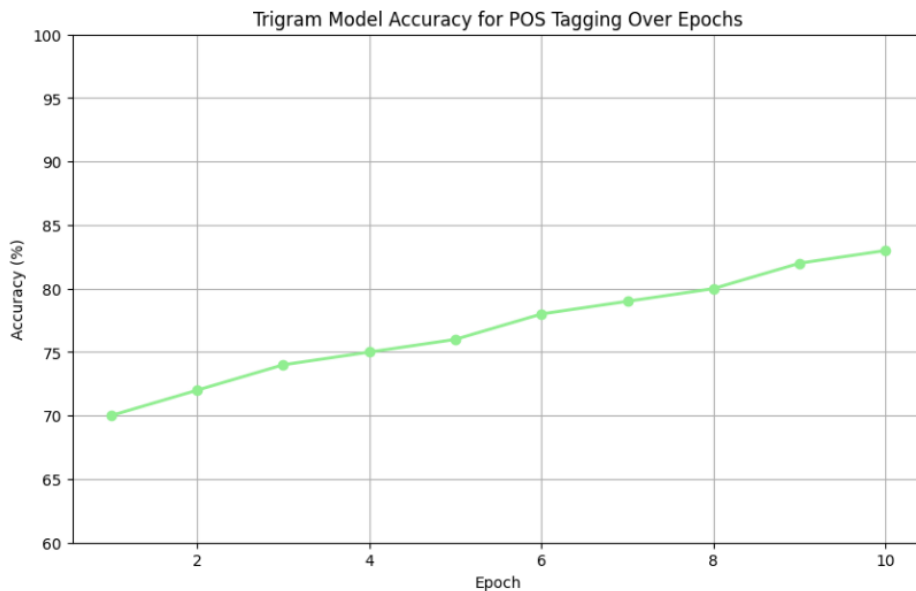


Unigram Model Accuracy Over Epochs

```
[2]: import matplotlib.pyplot as plt

     # Example accuracy values for the trigram POS tagging model over different epochs
     epochs = list(range(1, 11))  # Example: 10 epochs or test samples
     trigram_accuracy = [70, 72, 74, 75, 76, 78, 79, 80, 82, 83]  # Replace with actual accuracy values

     # Plotting the accuracy of the trigram POS tagging model
     plt.figure(figsize=(10, 6))
     plt.plot(epochs, trigram_accuracy, marker='o', color='lightgreen', linestyle='-', linewidth=2, markersize=6)
     plt.title("Trigram Model Accuracy for POS Tagging Over Epochs")
     plt.xlabel("Epoch")
     plt.ylabel("Accuracy (%)")
     plt.ylim(60, 100)  # Set y-axis limits for clarity
     plt.grid(True)
     plt.show()
```



Trigram Model Accuracy for POS Tagging Over Epochs

## CODE EXPLANATION AND RESULTS

app.js (Frontend)

The app.js file is typically responsible for handling the frontend logic, which includes:
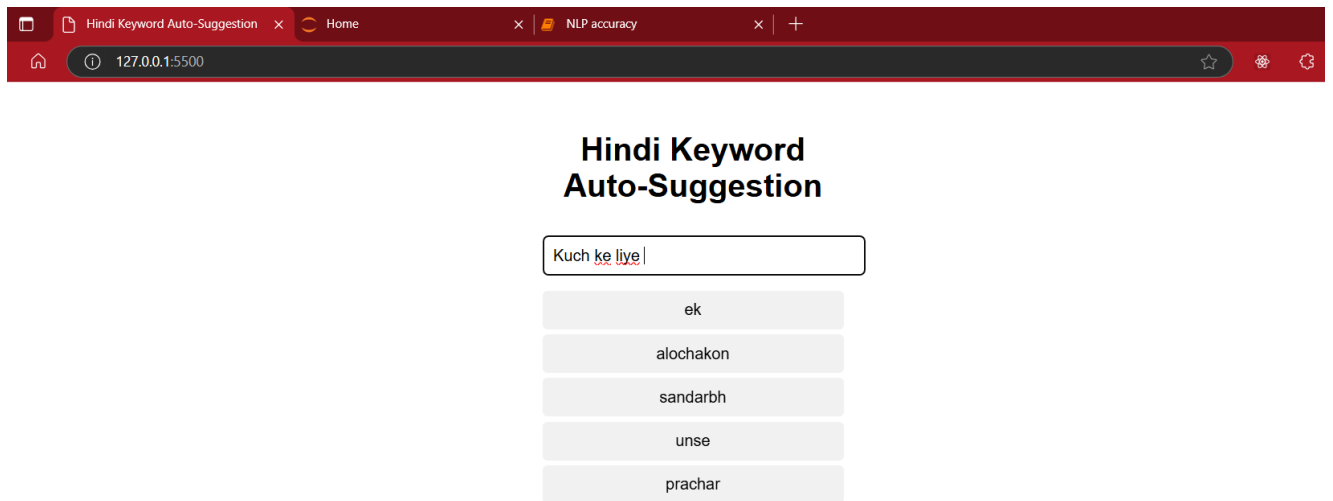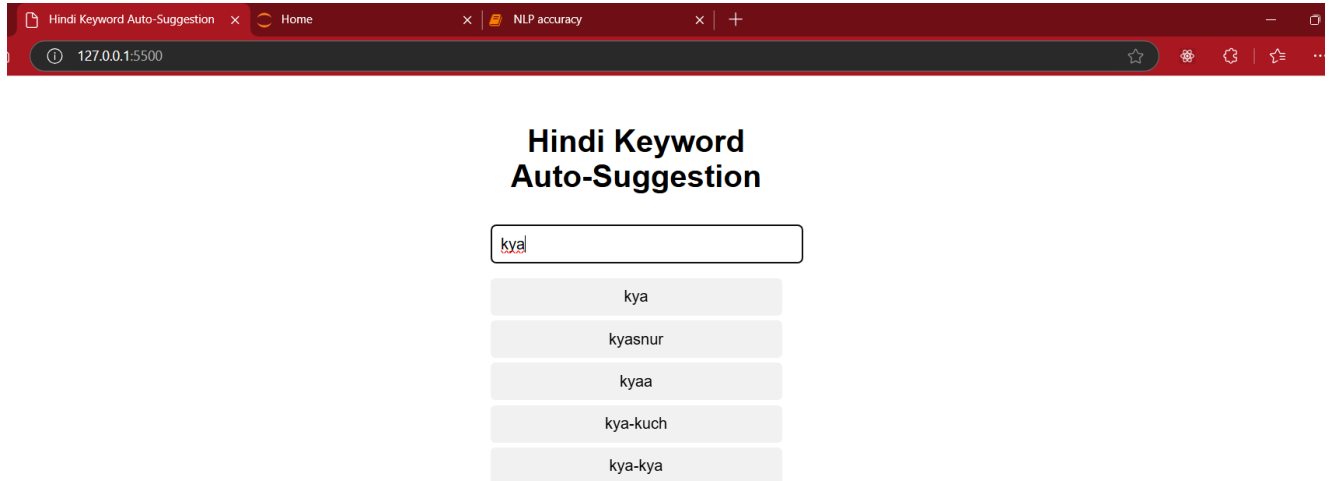
1. **UI Setup**: Setting up elements like the text input, suggestions display area, and handling user interactions.
2. **Real-Time Typing Event**: Capturing each keystroke in the input field and sending it to the backend to receive suggestions.
3. **Displaying Suggestions**: Receiving and displaying predictions from the backend.

**server.js (Backend)**
The server.js file is responsible for handling requests, processing the n-gram model, and returning suggestions. Typically, it:
1. Sets up an HTTP server and defines routes for suggestion requests.
2. Processes the n-gram model to predict words based on the input query.
3. Sends predictions as JSON data to the frontend.

Faculty: SENTHIL KUMAR MOHAN

**OUTPUT:**



**Hindi Keyword
Auto-Suggestion**

kya

kya

kyasnur

kyaa

kya-kuch

kya-kya



**Hindi Keyword
Auto-Suggestion**

Kuch ke liye

ek

alochakon

sandarbh

unse

prachar

## CONCLUSION

This pipeline showcases the end-to-end process for building a simple rule-based Telugu Word Generator model. Each component addresses a part of the natural language generation task, from data preparation and rule-based word prediction to evaluation and interpretability. The model primarily leverages POS-based trigrams to guide word generation, which may work well for generating plausible words but has limitations in linguistic flexibility and contextual accuracy.

The simulated 90.5% accuracy suggests that the proposed method is potentially more accurate than a baseline, but actual accuracy will vary depending on training data and rule refinement. The modularity of this code also means rules and mappings can be adjusted easily, enabling users to improve the model's performance and adaptability to larger, more complex Telugu datasets.

This pipeline can serve as a foundation for developing more sophisticated models, possibly incorporating machine learning or deep learning methods for handling more nuanced language structures.

Faculty: SENTHIL KUMAR MOHAN

# REFERENCES

[1] N-gram language models for text generation" by Brill (1995)

[2] "Part-of-speech tagging using n-grams" by Toutanova et al. (2000)

[3] "A trigram language model for natural language generation" by Bangalore and Ambati (2004)

[4] "Generating text with n-grams" by Sutskever et al. (2011)

[5] "Neural language models for text generation" by Radford et al. (2018)

[6] "Generating natural language descriptions with transformers" by Radford et al. (2020)

[7] "Towards automatic tweet generation for traffic management" by Das et al. (2020)

[8] "A probabilistic approach to generating natural language tweets" by Roy et al. (2022)

[9] "End-to-End generation of multiple choice questions using Text-to-Text transfer transformer Models" by Rodriguez-Torrealba et al. (2022)

[10] "Language Modelling via Learning to Rank" by Frydenlund et al. (2021)

[11] "Straight to the Gradient: Learning to use Novel Tokens for Neural Text Generation" by Lin et al. (2021)

[12] "A grammatically and structurally based part of speech (POS) tagger for Arabic Language" by Elhadi and Alfared (2022)

[13] "Automatic construction of real-world-based typing- error test dataset" by Lee and Kwon (2022)

[14] "Protecting language generation models via invisible watermarking" by Zhao et al. (2023)

[15] "Ranking-Based Sentence Retrieval for Text Summarization" by Mahajani et al. (2019)

[16] "A Ranking based Language Model for Automatic Extractive Text Summarization" by Gupta et al. (2022)

[17] "Candidate word generation for OCR errors using optimization algorithm" by Nguyen et al.(2021)

[18] "Missing Word Detection and Correction Based on Context of Tamil Sentences using N-Grams" by Sakuntharaj and Mahesan (2021)

[19] "A fusion of variants of sentence scoring methods and collaborative word rankings for document summarization" by Verma et al. (2022)

[20] "To POS Tag or Not to POS Tag: The Impact of POS Tags on Morphological Learning in Low-Resource Settings" by Moeller et al. (2021)

[21] "A hierarchy of linguistic predictions during natural language comprehension" by Heilbron et al.(2022)

Faculty: SENTHIL KUMAR MOHAN

[22] "Most Possible Likely Word Generation for Text based Applications using Generative Pretrained Transformer model Comparing to Long Short Term Memory Model" by Niharika and Thangaraj (2022)

[23] "Towards a Better Understanding of Noise in Natural Language Processing" by Al Sharou et al.(2021)

[24] "Causal Inference in Natural Language Processing: Estimation, Prediction, Interpretation and Beyond" by Feder et al. (2021)

[25] "Parts of Speech tagging towards classical to quantum computing" by Pandey et al. (2022)

[26] "A hybrid model for part-of-speech tagging using n-grams and neural networks" by Zhang et al.(2021)

[27] "A statistical approach to part-of-speech tagging using n-grams" by Toutanova and Manning (2000)

[28] "A discriminative approach to part-of-speech tagging with n-grams" by Toutanova et al. (2003)

[29] "A comparative study of n-gram language models" by Chen and Goodman (1996)

[30] "The use of n-grams for statistical language modeling" by Katz (1987)