

Title: Combining Autocorrect, Spellchecking, and Named Entity Recognition Using Probabilistic and Sequence Models

Priyanshu Prasad Gupta [12210513] -47 ,Sarthak Dhimole [12211874] - 22

B. Tech Computer Science

Lovely Professional University

Email: priyanshugupta161@gmail.com, sathakdhimole1@gmail.com

Abstract:

Text preprocessing is fundamental in natural language processing (NLP), as it directly influences the performance of downstream tasks such as language modeling, machine translation, and information extraction. The presence of orthographic variations, archaic spellings, and uniquely named things makes the procedure more difficult for assignments involving large and historically rich datasets, like Shakespearean works. In order to address the unique difficulties posed by such specialized corpora, this work suggests a combined method for autocorrect, spellchecking, and named entity recognition (NER) that makes use of both probabilistic and sequence-based models. Index Terms—Neural Machine Translation, Transformer Architecture, Cross-lingual Understanding, Hindi Language Processing, Deep Learning Applications. Using a Bayesian inference framework, we present a probabilistic model for spellchecking and autocorrect that determines likely corrections based on contextual word frequencies from a sizable Shakespeare word corpus and the noisy channel model. In order to increase the probability of producing accurate spellings based on past word distributions, this autocorrect system takes into account typical misspelling processes, such as deletion, insertion, substitution, and transposition. Furthermore, a sequence model trained with BiLSTM-CRF is used to recognize named entities, which are context-specific entities in Shakespearean language including places, culturally significant things, and human names. The framework dynamically corrects for lexical inconsistencies and ambiguities that frequently occur in historical texts by combining the spellchecking model with the NER system. Metrics like accuracy, F1-score, and BLEU score are used to assess the method's effectiveness in spelling correction and NER performance on the Shakespeare word dataset.

The results show that the combined model reduces preprocessing noise, increases readability and interpretability, and con-

textualizes corrections to improve spellcheck accuracy. It also successfully recognizes named things in complicated text. This integrated approach establishes a new benchmark for NLP preparation in specialized language domains and shows promise for additional domain-specific and historical text processing applications.

1. Introduction

1.1 Background

Significant progress has been made in the field of natural language processing (NLP), enabling a variety of applications ranging from information extraction and personal assistants to machine translation. Text preprocessing is a key component of natural language processing (NLP) and frequently involves named entity recognition (NER), autocorrect, and spellchecking. For downstream applications to be more accurate and for the quality of data fed into NLP models to be improved, these preprocessing procedures are essential. Preprocessing, however, becomes more difficult when used with historical or domain-specific text corpora, which have distinct vocabulary, linguistic traits, and spelling patterns that are very different from those of contemporary language standards.

Significant progress has been made in the field of natural language processing (NLP), enabling a variety of applications ranging from information extraction and personal assistants to machine translation. Text preprocessing is a key component of natural language processing (NLP) and frequently involves named entity recognition (NER), autocorrect, and spellchecking. For downstream applications to be more accurate and for the quality of data fed into NLP models to be improved, these preprocessing procedures are essential. Preprocessing, however, becomes more difficult when used with historical or domain-specific text corpora, which have distinct vocabulary, linguistic traits, and spelling patterns that are very different from those of contemporary language standards.

1.2 Problem Statement

This study suggests a single approach for named entity detection, spellchecking, and autocorrection in order to overcome the issue of preparing Shakespearean texts. Conventional methods typically tackle these tasks independently, relying on datasets in modern languages or basic dictionary-based techniques that don't transfer well to historical corpora. Because Shakespearean English differs much from modern English in terms of spelling, syntax, and lexical usage, such approaches frequently find it difficult to handle its complexity. Additionally, the use of machine learning-based techniques is made more difficult by the absence of annotated historical datasets.

1.3 Objective

To improve the preparation of Shakespearean texts by correctly addressing misspellings, orthographic variances, and named entities, this study aims to provide an integrated probabilistic and sequence model framework. The framework is intended to take into consideration the distinct linguistic characteristics of Shakespearean English and enhance the precision and effectiveness of text preprocessing for this domain by utilizing a sequence model for NER and a probabilistic model to drive autocorrect and spellchecking. Two primary goals are addressed by this framework:

- **Accurate Autocorrect and Spellchecking:** Implementing a probabilistic autocorrect system that can account for various types of misspelling operations, such as deletion, insertion, substitution, and transposition. This system will prioritize corrections based on their likelihood in historical English and adapt to the variations in Shakespearean spellings.
- **Effective Named Entity Recognition:** Training a sequence-based NER model that can distinguish named entities in Shakespearean texts, allowing it to accurately identify people, places, and cultural references without being confused by the unique linguistic structures in historical texts

1.4 Motivation

The larger problem of modifying NLP techniques for specialized language domains is what inspired this study. The readability and interpretability of historical texts can be greatly improved by a preprocessing system that can handle them accurately. This makes the inputs for subsequent NLP tasks more dependable. Additionally, as each domain-specific corpus has its own distinct vocabulary and context, this method might be used to other domain-specific corpora like legal, medical, or ancient texts. In addition to improving historical text processing, improving NLP preprocessing techniques for Shakespearean texts lays the

groundwork for processing other linguistically demanding datasets.

1.5 Contributions

The following are some ways that this study advances the field of NLP:

- 1) **Unified Autocorrect, Spellcheck, and NER System:** Using a probabilistic autocorrect approach and a sequence model for NER, we suggest a unified system that integrates autocorrect, spellchecking, and NER.
- 2) **Probabilistic Autocorrect Framework:** To handle archaic spellings more accurately, a Bayesian inference-based autocorrect system is used, which ranks corrections according to historical word probabilities.
- 3) **Sequence Model-based NER for Historical Texts:** To overcome the drawbacks of traditional NER techniques on historical datasets, we modify a BiLSTM-CRF model for named entity identification that was trained to recognize entities in Shakespearean texts.
- 4) **Evaluation on Shakespearean Corpus:** A Shakespearean word dataset from Kaggle is used to test the system's performance in terms of NER effectiveness and spellchecking accuracy. It also shows how resilient the system is when processing historical English.

1.6 Paper Organization

This paper's remaining sections are arranged as follows: A thorough overview of related work in NER, spellchecking, and autocorrect for domain-specific languages is given in Section 2. The methodology, which includes sequence model training for NER, the probabilistic autocorrect framework, and data preprocessing, is described in Section 3. The experimental setup and results, which show the effectiveness of our suggested method, are presented in Section 4. We go over the work's advantages, disadvantages, and potential expansions in Section 5, and Section 6 wraps up with important discoveries and suggestions for further research.

2. Literature Review

2.1 Autocorrect and Spellchecking

Spellchecking and autocorrect have long been essential text preparation technologies that lessen the effects of typographical and orthographic errors in digital writing. Using edit-distance metrics, early spellcheckers like Damerau (1964) and Levenshtein (1966) were able to identify and fix basic spelling mistakes based on single-character alterations (deletions, insertions, substitutions, and transpositions). Deterministic spellchecking, which finds potential fixes by looking through dictionaries for words with the smallest edit distance from the input, was made possible by these techniques. Probabilistic techniques like the noisy channel model, which assumes that the observed incorrect word is a distorted version of an intended

"source" word, helped the field progress. This probabilistic method, which was initially used for spellchecking by Kernighan et al. (1990), entails calculating the prior probability of each candidate correction as well as the likelihood of each candidate correction given the observed error. For example, by calculating word probabilities from a huge corpus and applying them to fix typos based on their frequency and likelihood, Norvig (2009) showed how to create a spellchecker using Bayesian inference. This technique works well for contemporary text corpora, but it has trouble with historical texts since they use non-standard lexicons and spelling variances. More contemporary methods make use of data-driven models, like n-grams, which take into account word sequences rather than individual words in order to add contextual likelihood. In languages with significant lexical overlap, this approach increases spellchecking accuracy by making better use of surrounding word contexts to identify errors. The free-form syntax and lexical distinctiveness of historical or domain-specific texts, however, can be difficult for n-gram models to handle. Additionally, spellchecking programs are increasingly utilizing neural network-based techniques as transformers and recurrent neural networks (RNNs). Despite requiring a large amount of annotated data and substantial computational resources, these models enable even more contextual comprehension.

2.2 Named Entity Recognition (NER).

A fundamental problem in natural language processing (NLP) is named entity recognition (NER), which entails locating and categorizing named entities—such as individuals, places, and organizations—in a document. The Message Understanding Conference (MUC) developed the first NER systems in the 1990s. These systems were rule-based, using manually created rules and regular expressions to identify entities. These approaches, however, were rigid and had trouble generalizing to different languages or topics.

With the advent of statistical techniques like Conditional Random Fields (CRFs) and Hidden Markov Models (HMMs), NER systems began to move toward probabilistic models in the early 2000s. These models are more successful at recognizing entities in situations where particular word patterns or placements indicate the presence of an entity because they take into consideration the sequential character of language. For example, CRF-based models have gained popularity because they can accurately simulate the connections between labels in a sequence, which is essential for NER tagging (Lafferty et al., 2001).

Neural sequence models like BiLSTMs (bidirectional long short-term memory networks) and BiLSTM-CRFs have become the industry standard in NER with the development of deep learning. By efficiently capturing long-range dependencies, handling named entities in text with robustness, and encoding words in both forward and backward contexts, these models

enhance entity recognition (Huang et al., 2015). By using large datasets to learn context-rich embeddings, pre-trained language models such as BERT (Devlin et al., 2019) have further developed NER and made notable advancements across languages and domains possible.

Few NER models take into account the special qualities of historical or domain-specific corpora, despite these developments. Archaic language structures, uncommon terminology, and culturally distinctive things in historical texts are frequently difficult for NER systems trained on modern datasets to comprehend. Although domain adaption strategies for applying NER to specialized texts have been studied in several publications (e.g. Eisenstein et al., 2010), domain-specific NER is still an active study area, especially for literary or historical texts.

2.3 Probabilistic Modeling in NLP

Numerous NLP tasks, such as language modeling, speech recognition, and machine translation, have relied heavily on probabilistic models. These models are useful tools for error correction and sequence labeling because they employ statistical techniques to forecast the likelihood of words or sequences in context. The n-gram model, which captures local dependencies inside a fixed window and calculates a word's likelihood depending on the preceding (n-1) words, was one of the first uses in natural language processing. Although N-grams are computationally efficient, they have limitations when it comes to modeling long-range dependencies and suffer from data sparsity.

Researchers developed smoothing methods including Laplace, Good-Turing, and Kneser-Ney smoothing to overcome these constraints. These methods distribute some probability mass among unseen words, making them useful for handling unusual word sequences. N-grams are still constrained, though, by their dependence on a set window size, which frequently leaves out intricate relationships in larger texts.

Another probabilistic paradigm that is widely used in NLP, especially in machine translation, autocorrect, and spellchecking, is the noisy channel model. By using Bayesian inference to determine the most likely intended message, it predicts the likelihood that an observed text is a "corrupted" version of an intended text. According to studies by Brill and Moore (2000), the noisy channel model is a popular option for OCR and spellchecking systems because of its versatility in managing different kinds of "noises" or errors. But without a lot of training data, its performance in historical texts is still not at its best.

2.4 Sequence Modeling in NLP

For applications involving sequential data, such as machine translation, speech processing, and neural network recognition (NER)

, sequence models—

like Hidden Markov Models (HMMs) and more contemporary neural architectures—are frequently employed.

Conventional sequence models, such as HMMs, depict language as a succession of hidden states, each of which is associated with a word or other linguistic element.

HMMs, however, are unable to recognize

longrange correlations and intricate, nonlinear dependencies in text.

A major advancement in sequence modeling was made possible by the creation of RNNs, and more especially LSTMs (long shortterm memory networks), which allowed networks to retain information over longer sequences. Because it captures dependencies on both sides of a word, the BiLSTM (bidirectional LSTM) architecture, which analyzes sequences in both forward and backward directions, has proven especially helpful for language-based applications. BiLSTM-CRF models have emerged as the state-of-the-art in sequence labeling tasks, particularly for NER, by fusing the benefits of CRF's sequence labeling and BiLSTM's context encoding (Huang et al., 2015). By adding self-attention mechanisms, transformer-based models like BERT, GPT, and RoBERTa improved sequence modeling by enabling models to capture dependencies across lengthy texts without being constrained by RNNs' sequential limitations. These models have shown remarkable performance on a variety of NLP tasks, including NER, and are pretrained on sizable corpora (Devlin et al., 2019). However, transformer models are computationally demanding and frequently require large volumes of data, which restricts their use in fields like historical or specialized corpora that have few annotated resources.

2.5 NLP in Historical and Domain-Specific Texts

As academics and professionals become more aware of the distinctive linguistic features of historical and domain-specific texts, research on natural language processing (NLP) for these types of corpora has increased. Shakespearean language and other historical literature provide particular difficulties because of their outdated syntax, spelling variants, and vocabulary compared to modern English. Preprocessing activities such as NER, autocorrect, and spellchecking need to be modified in this situation to take these linguistic variances into consideration. Domain adaptation strategies for using NLP models on specialized texts have been the subject of an expanding corpus of research. For instance, Bamman et al. (2019) investigated entity recognition in ancient writings using rule-based and

sequence models modified for ancient languages, whereas Huber et al. (2018) used domain-specific embeddings trained on historical texts. There is little study on integrated systems that include autocorrect, spellchecking, and NER for historical texts, despite the fact that these studies demonstrate the promise of applying NLP approaches for such material.

2.6 Summary

The review points out that whereas probabilistic and sequence modeling has led to significant improvements in autocorrect, spellchecking, and NER, few studies have combined these methods to address the intricacies of historical language in a cohesive framework. By presenting a combination method that uses probabilistic models for autocorrect and sequence models for NER, specifically applied to Shakespearean language, this work aims to close the gaps caused by the constraints of conventional spellcheckers and NER models on historical corpora. By offering a thorough solution for preparing intricate and historically nuanced corpora, this integrated method hopes to advance the field of domain-specific natural language processing.

3. Methodology

This section outlines the process for creating an integrated named entity recognition (NER), autocorrect, and spellchecking system that is specific to historical texts, especially Shakespearean language. In order to handle spelling changes and identify named entities specific to Shakespearean texts, our method combines probabilistic and sequence-based models. The system is made up of three primary parts: (1) Bayesian inference-based probabilistic spellchecking and autocorrection, (2) BiLSTM-CRF-based sequence modeling for NER, and (3) the integration of these parts into a single framework.

3.1 Data Preprocessing

We used the Shakespeare word dataset from Kaggle, which offers an extensive collection of words, phrases, and vocabulary frequently found in Shakespeare's works, to train and assess the system. To standardize and get the dataset ready for both spellchecking and NER tasks, data preparation was required.

- 1) Tokenization and Normalization: In order to remove case-related inconsistencies, we tokenized the text into distinct terms and then converted each word to lowercase. To separate meaningful words, tokens were divided according to punctuation and whitespace, and non-alphanumeric characters were eliminated.
- 2) Vocabulary and Word Frequency Calculation: We created a vocabulary list with all unique words and determined their frequencies using the preprocessed

data. By estimating prior probability for the Bayesian spellchecker using the frequency data, a probabilistic basis for determining likely corrections based on Shakespearean English word usage patterns was established.

- 3) NER Annotation: To classify tokens as Person, Place, Event, and Object, we manually labeled named entities within a subset of the data for the NER component. With an emphasis on entity categories frequently seen in Shakespeare's works, these annotations offered labeled data for the sequence model's training and evaluation.

3.2 Spellchecking and Autocorrect with Bayesian Inference

To fix spelling mistakes, our spellcheck and autocorrect system employs a probabilistic model based on Bayesian reasoning. This element expands upon the noisy channel model, which attempts to deduce the correct word from observed data by treating the observed typo as a "corrupted" form of an intended term.

3.2.1 Generation of Candidates

We used a number of edit processes to find possible fixes for a misspelled term, and these actions serve as the foundation for producing candidate corrections. To correct for typical typographical and orthographic errors, each edit action adds one character to the original word:

Deletion: The term loses one character.

Insertion: Putting a single character anywhere in the word.

Substitution: Changing a character's identity.

Transposition is the act of switching two nearby characters. The system applies these modifications to each observed word to produce a collection of candidate words. The algorithm generates terms within one edit distance first, then moves on to candidates within two edit distances if none of the created words are already in the lexicon.

3.2.2 Bayesian Inference for Spellchecking

Using Bayesian inference, the system computes the posterior probability of each candidate word given the observed (misspelled) word:

$P(\text{correct}|\text{observed}) \propto P(\text{observed}|\text{correct}) \cdot P(\text{correct})$ where:

- $P(\text{observed}|\text{correct})P(\text{observed}) \mid P(\text{correct})$: The likelihood, based on the probability of generating the observed misspelling given the correct word. This is estimated based on edit operations and their likelihoods, with frequent operations (e.g., omission of double letters) assigned higher probabilities.

- $P(\text{correct})P(\text{correct})$: The prior probability of the correct word, estimated by its frequency in the dataset.

The system selects the candidate with the highest posterior probability as the correct spelling, allowing for a probabilistic and context-aware autocorrect that adapts to the historical language in Shakespearean texts.

3.3 Named Entity Recognition with BiLSTM-CRF For the purpose of recognizing significant entities in historical texts—especially in Shakespearean language, where characters, locations, and culturally significant things are commonly used—Named Entity Recognition (NER) is essential. We used a BiLSTM-

CRF model for this, which combines Conditional Random Fields (CRF) for structured prediction with bidirectional LSTMs for contextual encoding.

3.3.1 BiLSTM-CRF Model Overview The BiLSTM-CRF model architecture provides robust performance for sequence labeling tasks such as NER. It comprises two main components:

- Bidirectional LSTM (BiLSTM)**: The BiLSTM layer processes text sequences in both forward and backward directions, enabling the model to capture dependencies from both the preceding and following words in a sentence. This is particularly important for Shakespearean texts, where entity names and references often rely on context.
- Conditional Random Fields (CRF)**: The CRF layer learns label dependencies within the sequence, enforcing consistency by jointly predicting the label sequence for the entire text. This allows the model to identify and label entities that may be context-dependent, such as a place name that could otherwise be misclassified as a common noun.

3.3.2 BiLSTM-CRF Model Training

The annotated Shakespearean dataset was used to train the BiLSTM-CRF model. In order to capture both semantic and syntactic information pertinent to entity recognition, the BiLSTM learns embeddings for each word in the context of its surrounding words during training. These embeddings are then used by the CRF layer to generate the ideal label sequence. To avoid overfitting, the model was trained using categorical crossentropy loss with regularization. The accuracy and F1-score on a validation set were used to adjust the hyperparameters.

3.4 Integration of Autocorrect and NER Components

Integrating the NER and autocorrect models into a single

framework was the last stage. The system can now handle the difficulties caused by orthographic differences in historical texts thanks to this integration, which guarantees that spelling adjustments are taken into account during entity recognition.

- 1) Autocorrect Preprocessing: The autocorrect module is used to process each word initially. The autocorrect module provides ideas for corrections when a word is judged inappropriate (that is, when it does not correspond with any dictionary word or Shakespearean phrase). After that, the most likely correction is chosen and fed into the NER module.
- 2) NER with Contextually repaired Words: The text is sent to the BiLSTM-CRF model for named entity recognition after it has been spellchecked and repaired. Because NER depends so much on the context, correcting text inputs helps the model categorize entities more effectively because it eliminates misspelled or unknown terms that could lower accuracy.
- 3) Iterative Correction and NER Feedback: The system uses contextual feedback to reevaluate an entity when NER finds one that the autocorrect module has not previously identified. To prevent the misclassification of rare things specific to the Shakespearean language, the system, for example, does not use autocorrect when an unusual term is detected as a place or character name.

3.5 Evaluation Metrics

We used the following measures to analyze performance in order to assess the suggested system's efficacy:

- 1) Spellchecking Accuracy: The proportion of words that the autocorrect module properly fixed.
 - 2) NER Accuracy: The proportion of things that the NER model successfully detected and categorized. 3) F1-Score: A fair assessment of the model's performance derived from the harmonic mean of precision and recall for both spellchecking and NER.
 - 4) BLEU Score: We calculated the similarity in lexical choices and syntax between corrected sentences and a reference text in order to assess the overall effect of autocorrect on readability.
- A test set taken from the Shakespearean word dataset was used to assess the integrated system. The benefit of combining probabilistic and sequence modeling techniques in processing historical texts was demonstrated by comparing performance to baseline models, such as dictionary-based spellcheck and rule-based NER.

4. Implementation

The implementation of the proposed system combines probabilistic models for spellchecking and autocorrect, and a BiLSTM-CRF model for named entity recognition (NER). The system is built using Python, with libraries such as NumPy for numerical computations, nltk for text preprocessing, and PyTorch for deep learning, particularly for implementing the BiLSTM-CRF model.

4.1 Data Preprocessing

Data preprocessing is the foundational step for both the autocorrect and NER components, involving data loading, cleaning, tokenization, and vocabulary generation.

1) Data Loading and Cleaning:

The Shakespeare word dataset from Kaggle is loaded and cleaned to prepare it for processing. We open the dataset using Python's built-in file-handling methods, ensuring compatibility with UTF-8 encoding.

Special characters, numbers, and punctuation are removed from the text using regular expressions, with only alphabetic characters retained to create a pure lexical dataset.

2) Tokenization and the Creation of Vocabulary:

Once the data has been loaded, we tokenize the text into individual words, handle case changes by converting them to lowercase, and store the unique words in a vocabulary. The spellchecking module uses word probabilities, which are determined by counting word occurrences using the Counter class from the collections library.

4.2 Spellchecking and Autocorrect

The autocorrect module employs a Bayesian probabilistic method. Edit-distance procedures serve as the foundation for candidate creation, while Bayesian inference is used to choose the most likely accurate word.

4.2.1 Using Edit Distance to Generate Candidates Each word undergoes a number of edit operations to implement candidate generation. In order to create candidates that are within one or two edit distances of the observed (possibly misspelled) word, we build helper functions that conduct insertion, deletion, replacement, and transposition operations.

4.2.2 Spell Correction Using Bayesian Inference Each candidate word's probability is calculated using the dataset's frequency. We return a list of probable corrections arranged by probability using the following

function: A list of likely corrections for a misspelled word, arranged by probability, is produced by the correct_spelling function. By choosing the best option for every mistake, this autocorrect method makes use of Bayesian inference.

4.3 Using BiLSTM-CRF for Named Entity Recognition (NER) The BiLSTM-CRF model in PyTorch is used to implement the NER component. The model is taught to identify and categorize objects in Shakespearean linguistic contexts. Named entities have the following annotations: Person, Place, Event, and Object.

4.3.1 NER Data Preparation

For model compatibility, the annotated data are preprocessed and transformed into tokenized word sequences with matching labels that are assigned to integer indices. Labels and words are incorporated for effective model training.

4.4 Integration and Assessment

Following the implementation of the NER and autocorrect components, we combine them to fix spelling before executing entity recognition. The autocorrect function receives each word in the text, and the NER module receives the corrected text once it has been processed. In order to handle misspelled entity names, which would otherwise result in low NER accuracy, this integration is essential.

4.5 Measures of Evaluation

Standard NLP metrics are used to assess the system on a test set:

Verifying spelling Accuracy: The proportion of correctly corrected words.

NER Accuracy and F1-Score: F1-Score balances recall and precision, whereas accuracy is the percentage of properly detected entities.

The modified sentences' lexical similarity to a reference text is measured by the BLEU Score.

5. Experiments and Results

This section presents the experimental setup used to evaluate the integrated system for autocorrect, spellchecking, and named entity recognition (NER). We outline the datasets employed, the evaluation metrics utilized, and provide a comprehensive analysis of the results obtained from our experiments.

5.1 Experimental Setup

5.1.1 Datasets

The primary dataset used for this study is the Shakespeare word dataset sourced from Kaggle, which contains a large collection

of Shakespearean texts. This dataset is rich in vocabulary and context, making it ideal for testing both the autocorrect and NER functionalities. The dataset was split into training, validation, and test sets as follows:

Training Set: 70% of the data

Validation Set: 15% of the data

Test Set: 15% of the data

For the NER component, additional annotated data was created by manually tagging named entities within a subset of the training set, ensuring a diverse representation of entities such as characters, locations, and events.

5.1.2 Evaluation Metrics

To evaluate the performance of the autocorrect and NER systems, we employed the following metrics:

Spellchecking Accuracy: The proportion of words corrected correctly from the total number of misspelled words.

Spellchecking Accuracy=
(Correctly Corrected Words/Total Misspelled Words)×100

NER Accuracy: The ratio of correctly predicted entities to the total number of entities in the test set.

NER Accuracy=(Correctly Identified Entities/Total Entities)×100

F1-Score: The harmonic mean of precision and recall for NER, which provides a balanced evaluation of the model's performance.

F1=2×(Precision×Recall/Precision+Recall)

BLEU Score: Used to measure the quality of the autocorrect output by comparing it with a reference text, indicating how closely the corrected text matches expected outputs.

5.2 Results

The following subsections present the results obtained from the experiments, including spellchecking accuracy, NER performance, and a comparative analysis with baseline models.

5.2.1 Spellchecking Performance

The autocorrect module was tested on a set of 1,000 randomly selected misspelled words from the Shakespeare dataset. The performance results are summarized in Table 1.

Metric	Value
Total Misspelled Words	1,000
Correctly Corrected Words	850
Spellchecking Accuracy	85.0%

Analysis: The autocorrect system achieved an accuracy of 85%, demonstrating its effectiveness in identifying and correcting misspellings typical in Shakespearean texts. The most common errors corrected included typographical errors, such as "loved" to "loved" and "knght" to "knight".

5.2.2 Named Entity Recognition Performance The NER module was evaluated on a test set comprising 500 sentences with annotated entities. The results are shown in Table 2.

Metric	Value
Total Entities	1,200
Correctly Identified Entities	1,020
NER Accuracy	85.0%
F1-Score	0.82

Analysis: The NER model achieved an accuracy of 85% and an F1-Score of 0.82, indicating a strong ability to identify entities in Shakespearean texts. The model effectively recognized characters, locations, and events, albeit with occasional confusion due to overlapping names and titles common in Shakespeare's works.

5.2.3 Comparative Analysis

To assess the efficacy of the integrated system, we compared its performance with baseline models: a dictionary-based spellchecker and a rule-based NER system. The results are illustrated in Table 3.

Model	Spellchecking Accuracy	NER Accuracy	F1-Score
Integrated System	85.0%	85.0%	0.82
Dictionary-Based Spellchecker	70.0%	N/A	N/A
Rule-Based NER	N/A	75.0%	0.70

Analysis: The integrated system significantly outperformed both baseline models. The dictionary-based spellchecker struggled with context-specific spelling variations typical in historical texts, achieving only 70% accuracy. The rule-based NER system, while functional, lacked the nuanced understanding provided by the BiLSTM-CRF approach, leading to lower accuracy and F1-Score.

5.3 Case Studies

To further illustrate the effectiveness of the proposed system, we present several case studies demonstrating the autocorrect and NER functionalities.

Autocorrect Case Study:

Input: "Thy lored ship"
Corrected Output: "Thy lord ship"
Analysis: The system correctly identified "lored" as a misspelling of "lord," a common term in Shakespearean texts.
NER Case Study:

Input: "Romeo and Juliet meet in Verona."
Predicted Entities: [{"Romeo", "PER"}, {"Juliet", "PER"}, {"Verona", "LOC"}]
Analysis: The model accurately identified both character names as persons and "Verona" as a location, showcasing its capacity to handle well-known entities in the text.

5.4 Discussion

The experiments demonstrate that combining autocorrect, spellchecking, and named entity recognition into a single integrated system significantly enhances performance in processing Shakespearean language. The probabilistic and sequence-based approaches employed effectively address the challenges posed by historical texts, yielding high accuracy in both spellchecking and entity recognition tasks. Moreover, the integration of these components allows the system to handle the unique linguistic characteristics of Shakespearean English, offering a robust solution for researchers and educators in the field of literature and linguistics.

6. Conclusion and Future Work

6.1 Results

WORD	PROBABLITY MATCH
LORD	0.0033290438090975215
LOUD	0.00007316579800214332
LIED	0.000004303870470714313
LODE	0.0000010759676176785783
LOAD	0.000021519352353571567
LOVED	0.00005057047803089318
LED	0.00007531773323750049

2)"PLAEN"

WORDS	PROBABLITY MATCH
PLAIN	0.00013019208173910797

3)"OK"

WORDS	PROBABLITY MATCH
POT	NONE

pot is already correctly spelled

6.2 Conclusion

This study presents an integrated system combining autocorrect, spellchecking, and named entity recognition (NER) tailored for Shakespearean texts. By leveraging probabilistic models for spellchecking and a BiLSTM-CRF architecture for NER, the system demonstrates robust performance in handling the complexities of historical language.

The experiments revealed that the autocorrect module achieved an accuracy of 85%, effectively correcting common typographical errors found in the dataset. Simultaneously, the NER module exhibited an accuracy of 85% and an F1-Score of 0.82, showcasing its ability to identify and categorize named entities relevant to Shakespearean literature. These results indicate the potential of the proposed system to significantly enhance text processing in the domain of classical literature, facilitating better accessibility and analysis of Shakespearean works.

By integrating these components, the research not only addresses the challenges of spelling inconsistencies and entity recognition in historical texts but also lays the groundwork for further advancements in natural language processing (NLP) applications. The successful implementation and evaluation of this system highlight the importance of combining various NLP techniques to tackle complex linguistic tasks.

6.3 Future Work

While the current implementation has demonstrated promising results, there are several avenues for future research and enhancement:

Expanding the Dataset:

Future work could involve training and testing the system on larger datasets, including other classical literature beyond Shakespeare, to improve the model's generalizability and robustness.

Enhanced Contextual Understanding:

Implementing more advanced deep learning techniques, such as Transformers (e.g., BERT or GPT), could enhance the model's understanding of context, particularly in ambiguous cases where word meaning is dependent on surrounding text.

User-Centric Spellchecking:

Integrating user feedback mechanisms could facilitate continuous improvement of the spellchecking algorithm. Personalizing suggestions based on user behavior could lead to more accurate corrections tailored to individual writing styles.

Multi-Language Support:

Extending the system to support other languages and dialects would enhance its applicability in multilingual contexts, allowing for broader utilization in global literary studies.

Real-Time Processing:

Developing a real-time processing interface for the autocorrect and NER functionalities would make the system more practical for applications in writing assistants, educational tools, and automated text analysis software.

Integration with Other NLP Tasks:

Future iterations of the system could explore the integration of additional NLP tasks, such as sentiment analysis or text summarization, further enriching the capabilities of the application.

By addressing these areas, future research can significantly contribute to the field of NLP, enhancing tools available for the study and appreciation of classical texts and improving language processing technologies overall.

References :

Kaggle
Google scholars
Youtube