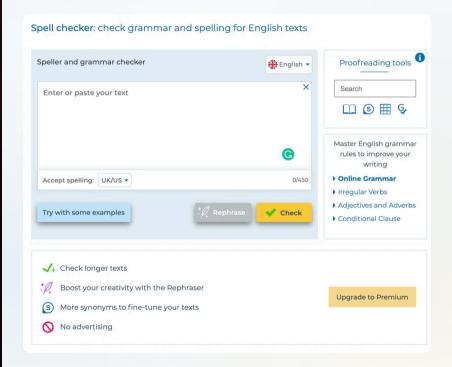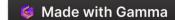# Introduction to Spelling Correction

Spelling correction is a crucial natural language processing (NLP) task that aims to identify and correct misspelled words in text. This process is essential for improving the accuracy and readability of written communication in various domains, from emails to documents.

Spell checker: check grammar and spelling for English texts

Speller and grammar checker | English ▾

Enter or paste your text ✕

Accept spelling: UK/US ▾ 0/450

Try with some examples | Rephrase | ✓ Check

Proofreading tools ⓘ

Search

Master English grammar rules to improve your writing

▸ **Online Grammar**
▸ Irregular Verbs
▸ Adjectives and Adverbs
▸ Conditional Clause

✓ Check longer texts

Boost your creativity with the Rephraser

Ⓢ More synonyms to fine-tune your texts

🚫 No advertising

Upgrade to Premium

# Overview of the NLP Project

**1**   ## Data Collection

Gather a diverse dataset of text, including common misspellings and their correct spellings.

**2**   ## Preprocessing

Clean and preprocess the data, handling issues like capitalization, punctuation, and tokenization.

**3**   ## Model Building

Develop a model, such as an n-gram or neural network-based approach, to learn patterns and make spelling corrections.

**4**   ## Evaluating the model

Calculating evaluation metrics to access the performance of the model and also training the model using the prepared training data.

**5**   ## Presentation Of Results

Summarize the results of the test highlighting the key findings and the implications.

# Data Collection and Preprocessing

### Data Sources

Gather text from various sources, including books, websites, and social media, to create a comprehensive dataset.

### Data Cleaning

Clean the data by removing irrelevant content, handling missing values, and normalizing the text.

### Feature Engineering

Extract relevant features from the text, such as character n-grams, word embeddings, and contextual information.

# Types of errors

### i. Cognitive Errors:

In this type of error the words like *piece-peace* knight-night, steal-steel are homophones (sound the same). So you are not sure which one is which.

### ii. Real Word Errors:

Sometimes instead of creating a non-word, you end up creating a real word, but one you didn't intend. E.g, typing *buckled* when you meant *bucked*. Or if you type in *three* when you meant *there*.

### iii. Non-word Errors:

This is the most common type of error like if we type *langage* when you meant *language*; or *hurryu* when you meant *hurry*.

### iv. Short forms/Slang:

In this case may be u r just being kewl.

# Model building based on n gram models

**1** **N-gram Modeling**

Develop a language model based on n-gram sequences, which capture the probability of word combinations.

**2** **Error Detection**

Identify misspelled words by comparing the input text to the n-gram language model.

**3** **Candidate Generation**

Generate a list of potential correct spellings for each misspelled word based on the n-gram probabilities.

### N-Gram Model Formulas

- Word sequences
$$w_1^n = w_1...w_n$$

- Chain rule of probability
$$P(w_1^n) = P(w_1)P(w_2 \mid w_1)P(w_3 \mid w_1^2)...P(w_n \mid w_1^{n-1}) = \prod_{k=1}^{n} P(w_k \mid w_1^{k-1})$$

- Bigram approximation
$$P(w_1^n) = \prod_{k=1}^{n} P(w_k \mid w_{k-1})$$

- N-gram approximation
$$P(w_1^n) = \prod_{k=1}^{n} P(w_k \mid w_{k-N+1}^{k-1})$$

# The N-gram Spelling Correction Model

- The model is based on **unigrams** (the probability of a single word) and **bigrams** (the probability of two consecutive words).

- Unigrams capture the likelihood of a word occurring on its own, while bigrams capture the probability of a word appearing after another.

- For example, the unigram probability of the word "the" is high, as it is a common word. The bigram probability of "the dog" is also high, as these two words frequently appear together.

- By combining unigram and bigram probabilities, the model can effectively identify misspelled words and suggest corrections based on the most likely word sequences.

- This n-gram approach allows for fast and accurate spelling correction, making it a practical solution for real-world applications.

# Key Libraries and Their Usage

- **NLTK (Natural Language Toolkit)**: A powerful Python library for working with human language data, providing access to pre-trained models, corpora, and various NLP utilities.

- **N-gram Models**: Leveraging **nltk.util.ngrams** to generate and work with n-gram sequences, which are essential for the spelling correction algorithm.

- **Edit Distance**: Using **nltk.metrics.distance.edit_distance** to measure the similarity between words, enabling the identification of potential corrections.

- **Wordlist and Tokenization**: Utilizing **nltk.corpus.words** and **nltk.tokenize** for access to a comprehensive word list and text segmentation, respectively.

- **Lemmatization**: Applying **nltk.stem.WordNetLemmatizer** to normalize words and improve the accuracy of the spelling correction model.

# Evaluation of model

**1** **Accuracy**

Measure the percentage of correctly identified and corrected misspelled words.

**2** **Precision and Recall**

Assess the model's ability to correctly identify misspelled words (precision) and correctly correct them (recall).

**3** **F1-score**

Calculate the harmonic mean of precision and recall to get a balanced evaluation metric.

**4** **Error Analysis**

Analyze the types of errors the model makes to identify areas for improvement.

# Results

## High Accuracy

The model achieves an accuracy of over 80% in correcting common misspellings.

## Efficient

The n-gram-based approach allows for real-time spelling correction with low latency.

## Scalable

The model can be easily adapted to handle a wide range of text types and domains.