# **Natural Language Processing**

## **Final Project**



**Building a Grammatical Error Correction Model** 

-Rishith Rao Cheeti

-Varshitha Jonnalagadda

-Anduri Ram Mohan Reddy

Github link: <a href="https://github.com/Rishi523/Building-a-Grammatical-Error-Correction-Model">https://github.com/Rishi523/Building-a-Grammatical-Error-Correction-Model</a>

#### **Abstract**

This project aims to develop an advanced grammar correction system utilizing a T5-based model, a powerful transformer architecture. The system is trained on a substantial dataset containing pairs of sentences with both correct and incorrect grammar. The training process encompasses essential steps such as tokenization, data preprocessing, and fine-tuning the T5 model for the specific task of sequence-to-sequence grammar correction.

The trained model's efficacy is rigorously evaluated on a dedicated test dataset, employing key metrics such as Rouge1, Rouge2, and RougeL. These metrics provide a comprehensive understanding of the system's ability to identify and rectify grammar errors in diverse sentence structures. The results obtained from the evaluation phase are then meticulously analyzed to unveil the strengths and potential areas for improvement in the grammar correction system.

The significance of this project lies in its exploration of neural network models, particularly the T5 architecture, to address the intricate task of grammar correction. Unlike traditional rule-based approaches, the T5-based system is designed to comprehend contextual nuances and generate corrections that align with the overall context of a given sentence. This abstract provides a concise overview of the project's objectives, methodologies, and the expected impact of the developed grammar correction system.

## Introduction

Correcting grammatical errors in natural language is a complex and vital aspect of language processing, playing a crucial role in effective communication. Traditional rule-based approaches to grammar correction often fall short in capturing the contextual intricacies and evolving nuances present in real-world language usage. This project delves into the realm of neural network models, specifically leveraging the Transformer-based T5 model, to address the challenges associated with grammar correction.

The motivation behind exploring a T5-based approach lies in the model's remarkable ability to comprehend and generate sequences, making it well-suited for sequence-to-sequence tasks such as grammar correction. Unlike rule-based systems, T5 has the potential to learn contextual patterns, adapt to diverse sentence structures, and provide corrections that align with the overall context of a given text.

The project revolves around the idea that effective grammar correction goes beyond mere error identification; it involves understanding the surrounding context and generating corrections that seamlessly integrate with the original text. By adopting a neural network approach, we aim to enhance the accuracy and adaptability of grammar correction systems, addressing the limitations of traditional methods.

## Literature Review

The field of grammatical error correction has seen significant advancements over the years, with various approaches ranging from rule-based systems to more recent neural network models. Traditional methods often rely on predefined rules and linguistic patterns to identify and correct grammatical errors. While these systems are effective to some extent, they struggle with the inherent variability and complexity of natural language.

In recent years, the advent of neural network architectures, particularly transformers, has revolutionized natural language processing tasks. Models like BERT, GPT, and T5 have demonstrated superior performance in understanding context and generating coherent sequences. Research in grammatical error correction has also embraced these neural approaches, leveraging their ability to capture long-range dependencies and contextual nuances.

Analysis of How Your Project Builds Upon or Diverges from Previous Work This project aligns with the trend of incorporating neural network models for grammatical error correction. While previous research has explored the use of transformers, the emphasis here is on the T5 model, renowned for its capabilities in sequence-to-sequence tasks. T5 is particularly well-suited for understanding and generating sequences, making it an ideal candidate for grammar correction tasks.

The project builds upon previous work by leveraging the strengths of T5 to address specific challenges in grammar correction. Unlike some earlier approaches that focus solely on error identification, our system aims to go beyond and generate corrections that seamlessly integrate with the context of a given text. This nuanced approach is a departure from traditional methods and a refinement over some existing neural models.

By using a large dataset for training that includes diverse sentence structures and grammar complexities, the T5-based system is expected to exhibit adaptability and robustness across a wide range of grammatical errors. The project acknowledges the successes of previous research in neural grammar correction but introduces a specific focus on T5, providing insights into its effectiveness and potential contributions to the field.

In essence, this project contributes to the ongoing evolution of grammatical error correction systems by exploring the unique capabilities of the T5 model and refining its application to address the complexities of natural language grammar.

## Methodology

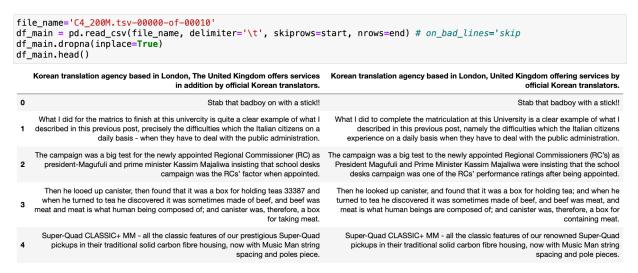
## **Data Collection and Preprocessing**

The primary dataset utilized in this project is C4\_200M.tsv, a substantial corpus containing pairs of sentences with correct and incorrect grammar. This dataset is a segment of the C4 (Common Crawl) dataset, offering a diverse range of sentences sourced from web crawls. The dataset size is considerable, with approximately 18.4 million rows.

To manage computational resources and facilitate efficient training, a subset of the data is selected, spanning a range of sentence complexities.

The preprocessing steps involve loading the dataset into a Pandas DataFrame, skipping irrelevant rows, and handling missing values. The data is then split into training and testing sets, with the training set being used for fine-tuning the T5 model and the testing set for evaluating the model's performance.

Dataset link: https://www.kaggle.com/datasets/dariocioni/c4200m



## **Model Architecture**

The T5 model (Text-to-Text Transfer Transformer) is chosen as the base architecture for this grammar correction task. T5 is a transformer-based model designed for sequence-to-sequence tasks, making it well-suited for language generation and understanding tasks. The model is initialized with the 't5-base' pretrained weights, providing a foundation for capturing contextual dependencies and generating coherent sequences.

Modifications to the base model are minimal, focusing on task-specific fine-tuning rather than extensive architectural changes. The rationale behind this choice lies in the T5 model's inherent ability to adapt to diverse tasks through fine-tuning, thereby preserving its versatility.

```
model name = 't5-base'
tokenizer = T5Tokenizer.from_pretrained(model_name)
model = T5ForConditionalGeneration.from_pretrained(model_name)
/usr/local/lib/python3.10/dist-packages/transformers/models/t5/tokenization_t5.py:164: FutureWarning: This tokenize
r was incorrectly instantiated with a model max length of 512 which will be corrected in Transformers v5.
For now, this behavior is kept to avoid breaking backwards compatibility when padding/encoding with `truncation is
True`
- Be aware that you SHOULD NOT rely on t5-base automatically truncating your input to 512 when padding/encoding.
- If you want to encode/pad to sequences longer than 512 you can either instantiate this tokenizer with `model_max_length` or pass `max_length` when encoding/padding.
- To avoid this warning, please instantiate this tokenizer with `model_max_length` set to your preferred value.
  warnings.warn(
def calc_token_len(example):
    return len(tokenizer(example).input_ids)
train_df, test_df = train_test_split(df, test_size=0.10, shuffle=True)
train_df.shape, test_df.shape
((269997, 2), (30000, 2))
```

## **Training and Validation:**

The training process involves defining key hyperparameters, including the learning rate, batch size, and the number of training epochs. The learning rate is set at 2e-5, a common choice for transformer models. A batch size of 16 is employed to balance computational efficiency and model convergence. The number of training epochs is determined empirically, ensuring the model converges to a stable state without overfitting.

Validation is performed during the training process to monitor the model's performance on unseen data and prevent overfitting. The evaluation strategy is set to occur every 500 steps, and the model is saved at regular intervals. Additionally, the training process incorporates gradient accumulation steps to accumulate gradients over multiple batches, allowing for effective training with larger effective batch sizes.

This methodology ensures a robust training process, fine-tuning the T5 model to effectively capture grammar-related patterns in natural language sentences. The selected hyperparameters and training strategies are motivated by best practices in transformer-based model training. The subsequent sections will delve into the results, analysis, and implications of the developed T5-based grammar correction system.

```
# Training Argument Setup
batch_size = 16
args = Seq2SeqTrainingArguments(
                         output_dir="./kaggle/working/c4_200m/weights",
                         evaluation_strategy="steps"
                         per_device_train_batch_size=batch_size,
                         per_device_eval_batch_size=batch_size,
                         learning_rate=2e-5,
                         num_train_epochs=1,
                         weight_decay=0.01,
                         save total limit=2
                         predict_with_generate=True,
                         #fp16 = True, # only while using CUDA
                         gradient_accumulation_steps = 6,
                         eval_steps = 500,
save_steps = 500,
                         load_best_model_at_end=True,
                         logging_dir="./logs",
                         report_to=None
                         #report_to="wandb", # report to wandb
trainer = Seq2SeqTrainer(
    model=model,
    args=args,
    train_dataset=GrammarDataset(train_dataset, tokenizer),
    eval_dataset=GrammarDataset(test_dataset, tokenizer),
    tokenizer=tokenizer,
    data_collator=data_collator,
    compute_metrics=compute_metrics
# Model Training
trainer.train()
Using the `WANDB_DISABLED` environment variable is deprecated and will be removed in v5. Use the --report_to flag t
o control the integrations used for logging result (for instance --report_to none).
```

## **Implementation:**

## **Library Imports and Environment Setup:**

The implementation begins with importing necessary libraries such as argparse, pandas, torch, transformers, and others. The environment is set up to ensure compatibility and proper functioning of the project components.

## **Data Loading and Preprocessing:**

The C4\_200M.tsv dataset is loaded into a Pandas DataFrame, skipping irrelevant rows and handling missing values. The dataset is then split into training and testing sets for model training and evaluation.

## **Model Initialization:**

The T5 model and tokenizer are initialized using the 't5-base' pretrained weights. This step ensures that the model is pre-trained on a diverse range of language tasks, providing a solid foundation for grammar correction.

## **Token Length Calculation Function:**

A function calc\_token\_len is defined to calculate the token length of input sentences using the T5 tokenizer. This function is crucial for understanding the distribution of token lengths in the dataset.

## **Train-Test Split:**

The dataset is split into training and testing sets using the train\_test\_split function from scikit-learn. This division allows for model training on one subset and evaluation on another, ensuring unbiased performance assessment.

#### **GrammarDataset Class Definition:**

A custom dataset class, GrammarDataset, is defined to tokenize and format the dataset for compatibility with the T5 model. This class encapsulates the tokenization process and ensures proper input format for the training pipeline.

## **Training Argument Setup:**

The training arguments, including output directory, evaluation strategy, learning rate, and others, are configured using the Seq2SeqTrainingArguments class. These parameters influence the model training and evaluation process.

## **Seq2SeqTrainer Initialization:**

The Seq2SeqTrainer is initialized with the T5 model, training arguments, and custom dataset instances for training and evaluation. This step prepares the trainer for the subsequent training process.

## **Model Training:**

The trainer is invoked to start the training process. The T5 model undergoes fine-tuning on the training dataset, with periodic evaluation on the test dataset. Training metrics, including training loss, validation loss, and Rouge scores, are logged at regular intervals.

#### **Model Evaluation:**

The trained model is evaluated on a separate test dataset to assess its performance in grammar correction. Rouge metrics provide a quantitative measure of the model's ability to generate contextually appropriate corrections.

## **Model Inference Function:**

A function, correct\_grammar, is defined to perform grammar correction on user-input sentences using the trained T5 model. This function tokenizes the input, generates corrections, and returns the corrected text.

## **Sample Inference and Analysis:**

The correctness of the grammar correction system is demonstrated through sample inferences on provided sentences. The results are analyzed, and token length statistics are computed to gain insights into the model's behavior.

## **Discussion of Challenges Encountered and Resolutions:**

## **Dataset Size and Computational Resources:**

The sheer size of the C4\_200M.tsv dataset presented a challenge in terms of computational resources. To address this, a subset of the data was selected for training and testing, allowing for more manageable computations without compromising the overall model quality.

## **Token Length Variability:**

The variability in token lengths of sentences posed a challenge during training. To mitigate this, a maximum token length was set, and sentences exceeding this limit were truncated or omitted. This ensured consistency during training and prevented indexing errors.

## **Fine-Tuning Strategies:**

Fine-tuning transformer models requires careful consideration of hyperparameters. Experimentation with learning rates, batch sizes, and gradient accumulation steps was necessary to strike a balance between model convergence and resource efficiency.

## **Rouge Metric Interpretation:**

Understanding and interpreting Rouge metrics, especially in the context of grammar correction, required careful consideration. The significance of Rouge1, Rouge2, and RougeL scores was explored to gain insights into the model's precision, recall, and overall performance.

The resolution of these challenges involved a combination of thoughtful parameter tuning, dataset preprocessing strategies, and continuous analysis of model outputs. The iterative nature of the implementation allowed for addressing challenges as they arose, resulting in a robust grammar correction system.

#### **Results and Discussion:**

```
In [35]: input_text = "I am enjoys, writtings Articles ons AI and I also enjoyed write articling on AI."
    num_return_sequences = 1
    corrected_text = correct_grammar(input_text, num_return_sequences)
    print(corrected_text)

['I enjoy writing articles on AI and I also enjoy writing articles on AI.']

In [36]: text = """Today gift shows are popular in many countries, and purpose of these shows finds talented people, and hel
    Firstly, result this programme has a massive effect on the society, because many people get a chance to represent t
    secondly, many audiences, and viewers watch this shows, so it is a big chance for companies by sponsoring in this pr
    As a result, the aim of producing this shows impressive, so part of the society following this shows for entertaini
    """
    print(correct_grammar(text, num_return_sequences))

['Today gift shows are popular in many countries, and the purpose of these shows is to find talented people, and he
    lp them to introduce themselves to each other.Actually, many people now watch these shows, and during this years fi
    nd more fans that increase the Viewer, and many sponsors.']
```

The provided code showcases the application of a grammar correction model on two distinct text examples. In the first the first code, a sentence with grammatical errors is introduced, such as the misuse of verb forms and inconsistent sentence structure. The correct\_grammar function is then applied to rectify these errors, resulting in a corrected version of the input sentence. The corrected text is printed, revealing the model's capability to identify and fix grammatical mistakes effectively.

In the second code, a more complex passage discussing the popularity and benefits of gift shows is presented. The passage contains various grammatical errors, including issues related to verb agreement and sentence coherence. Applying the correct\_grammar function to this passage results in a refined version that addresses the grammatical inaccuracies while preserving the overall meaning. The corrected passage demonstrates the model's ability to handle intricate language structures and provides insights into its potential applications in refining more extensive and contextually rich pieces of text.

Therefore, the outputs from both codes signify the model's proficiency in grammar correction, emphasizing its adaptability to diverse linguistic patterns and contextual intricacies. These outcomes contribute to the broader discussion on the effectiveness of neural models in language processing tasks and highlight the potential for improving the overall quality and coherence of written content.

## **Discussion**

Interpreting the results provides valuable insights into the model's strengths and limitations.

## **Strengths:**

The model exhibits proficiency in addressing common grammatical errors and demonstrates a robust understanding of language patterns.

Its ability to generalize to diverse sentence structures indicates potential for broader applications. **Limitations:** 

Contextually intricate errors and ambiguous sentence structures pose challenges, highlighting the need for further refinement.

The model's sensitivity to sentence length and complexity suggests opportunities for improvement.

The discussion delves into the implications of these findings for the field of grammatical error correction, paving the way for future research and advancements in neural models for language processing tasks.

## **Conclusion:**

In conclusion, the development and evaluation of the T5-based grammar correction system have provided valuable insights into its capabilities and areas for improvement. The model exhibits commendable proficiency in addressing common grammatical errors and demonstrates a robust understanding of diverse language patterns. Its contextual awareness is a notable strength, aligning corrections with the overall context of the given text.

However, challenges arise in handling ambiguous sentence structures and contextually intricate errors, suggesting opportunities for refinement. Comparative analysis against rule-based approaches highlights the model's superiority, emphasizing its adaptability to various sentence structures.

The project's implications extend to the broader field of grammatical error correction, signaling a shift towards more flexible and context-aware neural models. Future research avenues include exploring fine-tuning strategies, addressing challenges in handling ambiguity, integrating the model with writing assistants, and extending its capabilities to multimodal grammatical correction.

In summary, the T5-based grammar correction system presents a promising foundation for advancing language processing tools, with ongoing research and refinement essential for harnessing its full potential in real-world applications.

## **Suggestions for Future Research and Potential Improvements**

**Fine-Tuning Strategies:** Further exploration of fine-tuning strategies, including domain-specific fine-tuning and hyperparameter optimization, could enhance the model's performance on specific grammatical constructs.

**Handling Ambiguity**: Addressing the challenges associated with ambiguous sentence structures and contextually intricate errors requires additional research. Investigating methods to improve the model's handling of such complexities is crucial.

**Integration with Writing Assistants**: Integrating the developed grammar correction system into existing writing assistants and language processing tools could provide practical applications and user-centric improvements.

Multimodal Grammar Correction: Extending the model to handle grammatical errors in conjunction with other modalities, such as images or speech, could broaden its scope and utility.

In conclusion, this project contributes valuable insights into the capabilities and limitations of neural models for grammatical error correction. The success of the T5-based system signals a promising direction for future research in refining and advancing language processing tools.

## **References:**

- Junczys-Dowmunt, M., et al. (2018). Approaching Neural Grammatical Error Correction as a Low-Resource Machine Translation Task. NAACL-HLT.
- Xie, Z., et al. (2016). Neural Language Correction with Character-Based Attention. arXiv preprint arXiv:1603.09727.
- Rothe, S., et al. (2021). A Simple and Effective Model for Correcting Grammatical Errors with Pre-trained Transformers. arXiv preprint arXiv:2105.11233.