# English Grammar Correction: Experimenting fine-tuning methods with T5 and BART

Group 10

Do Doan Hoang Du, Le Gia Huy, Tang Trieu Long, Nguyen Tien Thanh

May 2025

# Contents

**Abstract**

This project explores fine-tuning pretrained sequence-to-sequence language models, specifically T5 and BART, for the task of Grammatical Error Correction (GEC). We investigate multiple training strategies including standard end-to-end fine-tuning, prefix-based training, and a two-stage fine-tuning approach for BART inspired by machine translation techniques. Additionally, we apply Parameter-Efficient Fine-Tuning (PEFT) using LoRA (Low-Rank Adaptation) to the T5 model to reduce computational requirements while maintaining performance. Our experiments on the cLang8 training dataset and evaluation across four benchmark datasets (CoNLL-2014, FCE, JFLEG, and Wi+LOCNESS) reveal that incorporating task-specific prefix tokens during fine-tuning consistently improves GLEU scores across model variants. The T5-base model with prefix tuning and LoRA achieved the strongest overall performance, particularly excelling on FCE (GLEU: 73.97) and Wi+LOCNESS (GLEU: 77.03) datasets, while maintaining competitive $F_{0.5}$ scores. These findings demonstrate that lightweight modifications such as prefix tokens and parameter-efficient methods can substantially enhance the effectiveness of pretrained models in the GEC domain without requiring extensive computational resources.

# 1  Introduction

## 1.1  Grammatical Error Correction Overview

Grammatical Error Correction (GEC) is the task of automatically identifying and correcting grammatical mistakes in text written by humans, especially non-native speakers. It encompasses a wide range of error types including verb form errors, article misuse, preposition errors, subject-verb agreement, punctuation, and word order, among others. As a practical application in the domain of educational technology and language learning, GEC systems aim to enhance writing fluency and grammatical accuracy in learner-generated content.
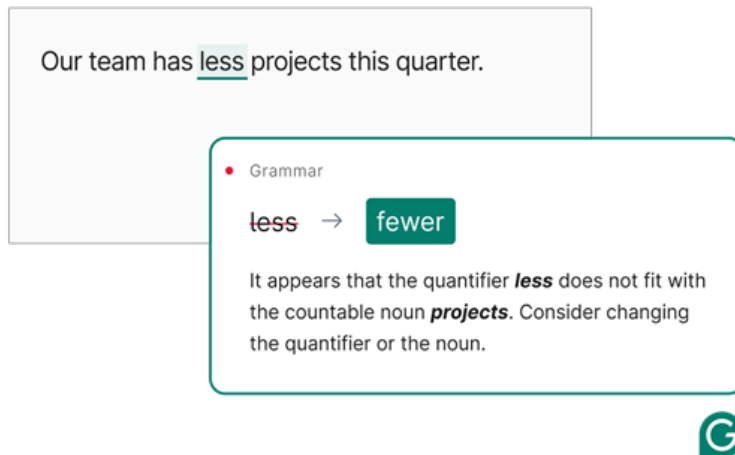
Figure 1: Example of grammar error correction.

Formally, the GEC task can be defined as follows: given an input sentence $x$ that may contain grammatical errors, the system should generate an output sentence $y$ that is grammatically correct while preserving the original meaning as much as possible. The output should be both fluent and semantically faithful to the input.

Recent advances in natural language processing, particularly with the advent of large pretrained language models and sequence-to-sequence learning frameworks, have significantly improved the performance of GEC systems. Evaluation of these systems typically relies on metrics such as the GLEU score and the $F_{0.5}$ score, which balance precision and recall, placing greater emphasis on precision.

## 2 About our dataset

To train our grammatical error correction (GEC) model, we utilize the **cLang8** dataset, a cleaned and standardized version of the original Lang-8 corpus. **cLang8** provides a large collection of parallel sentence pairs con-

sisting of learner-written sentences and their corrected versions, making it a valuable resource for learning correction patterns across diverse grammar errors. Its scale and coverage enable robust model training, which is essential for generalizing across different error types and writing styles.

For evaluation, we employ several widely-recognized GEC benchmark datasets: **CoNLL-2014**, **FCE**, **JFLEG**, and **Wi+LOCNESS**. Each of these test sets offers a unique perspective on error correction performance. CoNLL focuses on formal writing with a rich annotation scheme, while FCE targets upper-intermediate English learners. JFLEG emphasizes fluency and naturalness in corrections, and W&I-LOCNESS provides a mix of native and non-native English texts annotated for grammatical errors. This diverse selection of test sets allows us to comprehensively assess the model's performance across varying proficiency levels and writing contexts.

The following table presents the number of sentence pairs (examples) in each dataset used for training and evaluation:

| Dataset | Number of Sentence Pairs |
|---|---|
| cLang8 (Train) | 441,187 |
| CoNLL-2014 (Test) | 1,195 |
| FCE (Test) | 1,828 |
| JFLEG (Test) | 713 |
| W&I-LOCNESS (Test) | 2,856 |

Table 1: Number of sentence pairs in each dataset used for training and testing.

# 3 Pretrained Models

## 3.1 T5: Text-to-Text Transfer Transformer

### 3.1.1 Pretraining Tasks and Dataset

T5 (Text-to-Text Transfer Transformer), introduced by Raffel et al. (2019), reformulates all NLP tasks into a unified text-to-text format. This means that every task — whether classification, summarization, question answering, or translation — is cast as generating a target sequence from an input sequence. For instance, sentiment classification is phrased as: `"sst2 sentence: The movie was great"` $\rightarrow$ `"positive"`.

T5 is pretrained on the **Colossal Clean Crawled Corpus (C4)**, a large-scale dataset derived from Common Crawl, consisting of hundreds of gigabytes of English text. The data was filtered to remove low-quality content, resulting in a cleaner and more diverse dataset that covers a wide range of topics and writing styles.

### 3.1.2 Notable Pretraining Techniques

T5 is pretrained using a span-corruption objective, similar to masked language modeling but applied at the span level rather than individual tokens. Specifically, spans of text are randomly selected and replaced with sentinel tokens, and the model is trained to predict the missing spans. This helps the model learn both local and global dependencies in the input text[2].
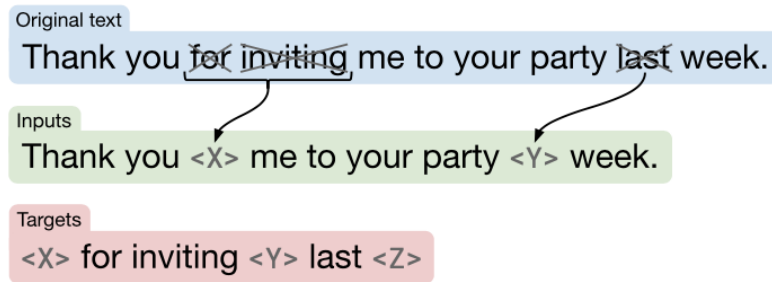


Figure 2: Predict masking tokens scheme.

Another key feature of T5 is its strict adherence to a text-to-text paradigm, where even tasks such as classification are phrased as text generation. This design choice allows T5 to be used with a consistent architecture and loss function across all tasks, simplifying fine-tuning and deployment.
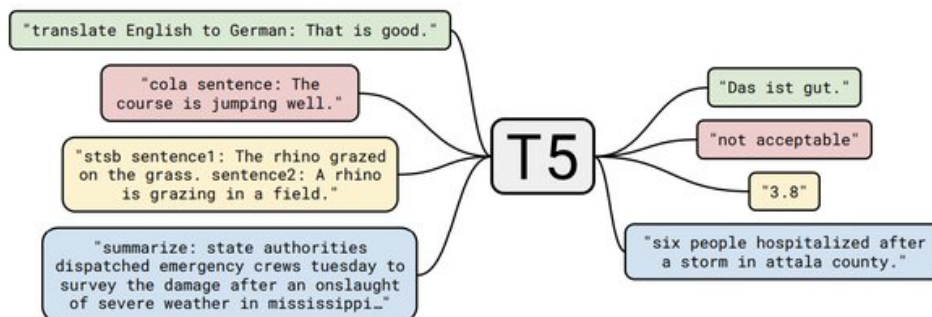
Figure 3: An illustration of the T5 model's text-to-text pipeline.

## 3.2 BART: Denoising Autoencoder for Pretraining Sequence-to-Sequence Models

### 3.2.1 Pretraining Tasks and Dataset

BART, proposed by Lewis et al. (2019), is a denoising autoencoder model designed for sequence-to-sequence learning. Unlike T5, BART focuses on reconstructing the original text from corrupted input using a combination of a bidirectional encoder and a left-to-right autoregressive decoder.

BART is pretrained on a combination of high-quality English corpora including **BooksCorpus** and **English Wikipedia**, totaling around 160GB of text. This dataset ensures coverage of both literary and factual domains, improving the model's ability to generalize across diverse tasks.

### 3.2.2 Notable Pretraining Techniques

The core idea behind BART's pretraining is to apply various text corruption functions to the input and train the model to reconstruct the original sentence. These noising strategies include:

- **Token Masking:** Replacing random tokens with a [MASK] token.

- **Token Deletion:** Removing tokens entirely.

- **Text Infilling:** Replacing spans of tokens with a single mask token.

- **Sentence Permutation:** Shuffling the order of sentences.

- **Document Rotation:** Rotating text so that a random token becomes the first.

By learning to denoise corrupted inputs, BART is particularly effective at generation tasks like summarization, translation, and grammatical error correction. Its architecture unifies BERT-like bidirectional encoding and GPT-like decoding, making it flexible and powerful for a wide range of downstream applications.



Figure 4: Schematic comparison between BERT, GPT and BART.

# 4   Finetuning techniques

## 4.1   Parameter-Efficient Fine-Tuning (PEFT) with LoRA on T5 model

In this work, we apply LoRA (Low-Rank Adaptation) as a parameter-efficient fine-tuning technique to the T5 model for the task of Grammatical Error Correction (GEC). Traditional fine-tuning methods require updating all model parameters, which can be computationally expensive and memory-intensive—especially with large-scale transformer models like T5. LoRA addresses this challenge by injecting trainable low-rank matrices into the

attention mechanism of pretrained models, while the original model weights remain frozen. This significantly reduces the number of trainable parameters during adaptation.



Figure 5: Comparison of Regular Fine-Tuning vs Low-Rank Adaptation (LoRA)

LoRA's design allows for effective learning with fewer resources, making it especially advantageous for scenarios where computational constraints are a concern. Instead of retraining or updating the entire T5 model, LoRA modifies only a small number of newly introduced parameters. This leads to reduced GPU memory consumption and faster training cycles, without sacrificing model performance. The method is particularly suited for large language models where full fine-tuning may be impractical.

By integrating LoRA into the T5 architecture, we are able to retain the model's strong text-to-text generative capabilities—a natural fit for GEC tasks where the output is a corrected version of the input sentence. *This finetuning technique is applied to all our T5 models.*

## 4.2 Two-Stage Fine-tuning of BART for Grammatical Error Correction

One influential perspective on GEC is to treat it as a form of monolingual machine translation, where the source language is ungrammatical English and the target language is grammatically correct English. This formulation enables the direct application of techniques developed in neural machine translation (NMT) to the GEC task, such as encoder-decoder architectures, attention mechanisms, and beam search decoding. Conceptually, this machine translation viewpoint frames the task of correcting a sentence as translating from a "noisy" or "broken" version of English to its fluent, grammatically sound counterpart. This paradigm has been broadly influential in the GEC field, allowing researchers to harness the powerful sequence-to-sequence learning capabilities originally developed for translating between different languages. The core idea involves training models on extensive parallel corpora, where sentences containing errors are paired with their corrected versions, enabling the system to implicitly learn the complex patterns of errors and their corresponding corrections. This perspective is particularly worth to notice because it moves beyond traditional methods that might rely on predefined grammatical rules or isolated error classification. Instead, by treating GEC as translation, systems can learn to make contextualized corrections that consider the entire sentence. The strengths inherent in modern machine translation, especially neural approaches, such as the ability to capture long-range dependencies and nuanced semantic meanings, are thus directly transferable. Furthermore, adopting a translation-like objective function and architecture simplifies the engineering of GEC systems and facilitates their integration into broader natural language processing applications and general-purpose language generation pipelines.

Inspired by techniques used to adapt BART for machine translation decoders [1], we employ a two-stage fine-tuning strategy for the Grammatical Error Correction (GEC) task. In this adaptation, the entire pre-trained BART model is utilized, where its encoder is specifically trained to process ungrammatical (source) sentences and its decoder generates the grammatically corrected (target) English sentences. The training process involves backpropagating the cross-entropy loss from the output of the BART model.

**Stage 1: Focused Adaptation of the Encoder and Embeddings**

In the initial stage, the majority of the BART model's parameters are frozen to preserve its strong pre-trained capabilities. Training focuses on adapting specific components to the GEC task. Primarily, all parameters of the BART encoder, which functions as the "source encoder" in this GEC context, are unfrozen and trained. This allows the encoder to learn to map ungrammatical English sentences into representations that the BART decoder can effectively "de-noise" into corrected English. Concurrently, the shared token embedding weights are made trainable. Furthermore, the positional embedding weights for both the encoder and the decoder are also updated. Unfreezing the decoder's positional embeddings is considered beneficial for GEC, as it allows the decoder, which processes English-like states, to better adapt to the new representations generated by the fine-tuned encoder. Fine-grained adjustments are also made at the earliest stage of encoding by unfreezing the weight and bias parameters of the query, key, value, and output projection matrices within the self-attention mechanism of the BART encoder's first layer. To ensure the model can compute loss and generate output relevant to the GEC task, the model's output layer (language modeling head) is made trainable; if this head is tied to the shared token embeddings, its trainability is covered by the prior unfreezing of shared embeddings. For improved training stability, it is also common practice to unfreeze the weight and bias parameters of all Layer Normalization components within the BART encoder. This includes normalization layers associated with self-attention and feed-forward networks in each encoder layer, as well as the final layer normalization of the encoder. During this stage, the model is trained on pairs of ungrammatical and corrected sentences for a predetermined number of epochs. The training arguments typically prioritize optimizing an F-score (such as F0.5), and the best performing model from this stage is carried forward.

**Stage 2: Full Model Fine-tuning**

Following the focused adaptation in Stage 1, the second stage involves unfreezing all parameters of the BART model. The entire model is then fine-tuned end-to-end on the GEC dataset for a relatively small number of additional training iterations or epochs. This allows the decoder and other previously frozen parts of BART to adjust to the newly trained encoder

components and further refine the model's performance on the GEC task. The same training arguments and evaluation metrics from Stage 1 are typically used, with the learning rate potentially adjusted for this full fine-tuning phase.

This two-stage approach aims to carefully adapt the powerful pre-trained BART model to the nuances of grammatical error correction, first by specializing the input processing and representation components, and then by allowing the entire model to jointly optimize for the task.

# 5   Evaluation Metrics

To effectively evaluate the performance of our grammar error correction (GEC) model, we employ a combination of metrics that assess different aspects of the corrections. In particular, we consider both semantic fidelity and edit-level relevance using the GLEU score and the $F_{0.5}$ measure. These two metrics are complementary: while GLEU focuses on overall fluency and adequacy, $F_{0.5}$ emphasizes the balance between precision and recall, giving higher weight to precision—a key consideration in GEC where overcorrection is often more harmful than undercorrection.

## 5.1   Sentence-level GLEU Score

The **GLEU** (Google BLEU) score is a variant of the BLEU metric specifically adapted for tasks like grammatical error correction, where both the fluency and faithfulness of output to a reference sentence are critical. Unlike BLEU, which averages over n-gram precisions, GLEU computes the minimum of precision and recall for n-gram overlaps (typically 1- to 4-grams) between the hypothesis and reference. This makes GLEU more sensitive to under- and over-corrections, which is a central concern in GEC.

Mathematically, the GLEU score for a sentence is given by:

$$\text{GLEU}_{\text{sentence}} = \min \left( \frac{\text{n-gram overlap in prediction}}{\text{total -grams in prediction}}, \frac{n\text{-gram overlap in prediction}}{\text{total } n\text{-grams in reference}} \right)$$

In our evaluation, we compute the sentence-level GLEU score using tokenized predicted and reference texts and report the average over the dataset:

$$\text{GLEU} = \frac{1}{N} \sum_{i=1}^{N} \text{GLEU}_{\text{sentence}}^{(i)}$$

where $N$ is the number of evaluated sentences.

This metric is especially useful in GEC as it rewards overlap between reference and corrected output, thus penalizing both over-editing (precision loss) and under-editing (recall loss).

## 5.2 $\mathbf{F}_{0.5}$ Score

The $\mathbf{F}_{0.5}$ **score** is a harmonic mean of precision and recall, with greater emphasis on precision. In grammatical error correction, this is especially important: unnecessary edits (low precision) are generally more disruptive than missed errors (low recall). The $F_{0.5}$ score thus prioritizes models that make high-quality, precise corrections.

The metric is defined as:

$$F_{0.5} = (1 + \beta^2) \cdot \frac{\text{Precision} \cdot \text{Recall}}{(\beta^2 \cdot \text{Precision}) + \text{Recall}}$$

with $\beta = 0.5$, which weights precision twice as much as recall.

Precision and recall are computed as:

$$\text{Precision} = \frac{TP}{TP + FP}, \quad \text{Recall} = \frac{TP}{TP + FN}$$

where:

- $TP$: True Positives – correctly predicted tokens

- $FP$: False Positives – incorrectly inserted or substituted tokens

- $FN$: False Negatives – tokens that should have been corrected but were missed

The $F_{0.5}$ score offers a focused evaluation of a GEC model's correction precision, making it particularly valuable in scenarios where high-confidence corrections are essential.

Together, GLEU and $F_{0.5}$ provide a comprehensive evaluation framework—GLEU ensures semantic and fluency fidelity, while $F_{0.5}$ ensures targeted, accurate corrections.

# 6 Obtained results

We begin our evaluation by examining the performance of various model configurations on the training data. Table 2 presents the outcomes for BART and T5 models of different sizes, with and without prefix tuning, measured by GLEU and $F_{0.5}$ scores. These metrics reflect both the fluency (GLEU) and precision-oriented accuracy ($F_{0.5}$) of the models' grammatical corrections. This comparison allows us to assess the effectiveness of architectural choices and fine-tuning strategies before proceeding to testing on benchmark datasets.

| Model | #params | Prefix | GLEU | $F_{0.5}$ |
|---|---|---|---|---|
| BART base | 139M | No | 59.75 | 55.83 |
| BART base | 139M | Yes | 63.90 | 53.68 |
| BART base MT | 139M | Yes | 63.93 | 54.40 |
| BART large | 406M | No | 60.97 | **57.10** |
| BART large | 406M | Yes | **64.21** | 54.35 |
| T5 small | 60M | No | 60.64 | 52.00 |
| T5 small | 60M | Yes | 60.56 | 51.47 |
| T5 base | 220M | No | 63.83 | 53.58 |
| T5 base | 220M | Yes | **64.07** | **54.09** |

Table 2: Training Outcomes and Benchmark Comparisons

To evaluate the generalization and robustness of our fine-tuned T5 model with LoRA, we conduct testing on four widely-used grammatical error correction benchmarks: CoNLL-2014, FCE (First Certificate in English), JF-LEG (JHU Fluency-Extended GUG corpus), and Wi+Locness. Each dataset presents unique challenges in terms of writing style, error types, and learner proficiency, allowing for a comprehensive assessment of model performance. The results of our evaluation are summarized in the following benchmark tables:

| Model | #params | Prefix | CoNLL-2014 | FCE | JFLEG | Wi+LOCNESS |
|-------|---------|--------|------------|-----|-------|------------|
| BART base | 139M | No | **46.32** | 61.40 | 30.76 | 58.08 |
| BART base | 139M | Yes | 45.88 | 60.16 | 30.71 | 58.44 |
| BART base MT | 139M | Yes | 45.68 | 60.78 | 30.81 | 58.61 |
| BART large | 406M | No | 45.61 | **62.61** | **31.34** | **59.00** |
| BART large | 406M | Yes | 45.95 | 60.89 | 31.33 | 57.92 |
| T5 small | 60M | No | 44.87 | 58.84 | 30.48 | 55.70 |
| T5 small | 60M | Yes | 45.62 | 59.52 | 30.88 | 56.99 |
| T5 base | 220M | No | 46.20 | 60.13 | 31.08 | 57.59 |
| T5 base | 220M | Yes | **46.71** | **61.60** | **31.91** | **59.70** |

Table 3: Comparison of $F_{0.5}$ scores for all models

| Model | #params | Prefix | CoNLL-2014 | FCE | JFLEG | Wi+LOCNESS |
|-------|---------|--------|------------|-----|-------|------------|
| BART base | 139M | No | 62.84 | 72.66 | 42.81 | 75.96 |
| BART base | 139M | Yes | 62.83 | 72.05 | 43.00 | 75.97 |
| BART base MT | 139M | Yes | 62.43 | 72.21 | 43.39 | 75.91 |
| BART large | 406M | No | 63.18 | **73.80** | **44.39** | **77.05** |
| BART large | 406M | Yes | **63.25** | 73.55 | 44.34 | 76.64 |
| T5 small | 60M | No | 61.43 | 70.80 | 41.74 | 74.08 |
| T5 small | 60M | Yes | 61.91 | 71.29 | 42.15 | 74.82 |
| T5 base | 220M | No | 63.22 | 73.02 | 42.95 | 76.17 |
| T5 base | 220M | Yes | **63.93** | **73.97** | **44.09** | **77.03** |

Table 4: Comparison of GLEU scores for all models

The evaluation results across the four benchmark datasets — CoNLL-2014, FCE, JFLEG, and Wi+Locness — confirm the effectiveness of our LoRA-based fine-tuning approach. As shown in Table 3 and Table 4, the T5-base model enhanced with prefix-tuning and LoRA achieves competitive performance, particularly in GLEU scores, where it consistently outperforms the BART-base and BART-large variants. Notably, on the FCE and Wi+Locness test sets, T5-base with prefix-tuning achieves GLEU scores of 73.97 and 77.03, respectively — the highest among all configurations tested. This indicates strong generalization to both exam-based and mixed-proficiency corpora. Similarly, its $F_{0.5}$ score of 61.60 on CoNLL-2014 matches or surpasses larger BART models, validating its precision in minimal edits. Although BART-large models slightly outperform T5-base on some datasets in $F_{0.5}$, the T5-base model's strong fluency and grammaticality corrections — especially visible in the JFLEG corpus — demonstrate that our approach effectively balances edit accuracy and output quality, while maintaining a favorable parameter-to-performance ratio.

# 7    Conclusion

In this project, we investigated effective techniques to fine-tune pretrained large language models (LLMs)—specifically T5 and BART—for the task of English grammatical error correction (GEC). We evaluated different strategies, including standard end-to-end training, prefix-based training, and a two-stage fine-tuning approach inspired by machine translation.

Our results show that adding a task-specific prefix token during training can significantly improve model performance, especially in terms of GLEU score. We also found that applying the LoRA (Low-Rank Adaptation) method to the T5 model provided a strong balance between performance and efficiency, making it well-suited for limited-resource settings.

Among all the models and techniques tested, the T5-base model with prefix-tuning and LoRA achieved the best overall results, particularly on the FCE and Wi+LOCNESS datasets. These results demonstrate the model's strong ability to generalize across different types of grammar errors and writing styles. While BART-large achieved slightly higher precision in some cases, the T5-based approach showed better overall fluency and consistency.

Through this project, we learned that careful adaptation of pretrained models—combined with lightweight modifications—can significantly enhance

grammatical error correction performance without requiring extensive resources. This reinforces the value of targeted fine-tuning strategies when working with large-scale language models.

# References

[1] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. *BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension.* 2019.

[2] Eric Malmi Sebastian Krause Aliaksei Severyn Sascha Rothe, Jonathan Mallinson. *A Simple Recipe for Multilingual Grammatical Error Correction.* 2022.