

# **TRƯỜNG ĐẠI HỌC THỦY LỢI**



HỌ VÀ TÊN: TRẦN THẾ MINH  
MÃ SINH VIÊN: 2251262619  
LỚP: 64TTNT2

## **BÁO CÁO BÀI TẬP LỚN HỌC PHẦN XỬ LÝ NGÔN NGỮ TỰ NHIÊN**

**SỬA CHỮA LỖI NGỮ PHÁP TRONG NGÔN NGỮ TIẾNG ANH**

NGƯỜI HƯỚNG DẪN: TS. Nguyễn Quang Hoan

Hà Nội, 2025

## LỜI NÓI ĐẦU

Trong thời đại xã hội ngày nay, ngôn ngữ nói chung và tiếng Anh nói riêng đã thực sự rất phổ biến và là ngôn ngữ ai cũng cần phải biết. Tiếng Anh đóng một vai trò cực kì quan trọng trong việc trao đổi văn hóa, giá trị, đời sống của tất cả mọi người trên thế giới. Nhưng đối với nhiều người vẫn còn sử dụng sai tiếng Anh đặc biệt là ngữ pháp, vậy nên em đã nhìn ra vấn đề này và viết mô hình sửa chữa lỗi sai ngữ pháp Tiếng Anh để giúp tất cả mọi người có thể sử dụng tiếng Anh một cách trơn tru và mượt mà nhất. Trong các mô hình như LSTM, Informer, Linear Regression,..... Em nhận thấy mô hình T5 (text-to-text transfer transformer) dựa trên kiến trúc Transformer. Mô hình này hoạt động rất tốt với việc đổi từ ngôn ngữ sang dạng chuỗi để có thể áp dụng ra mô hình như dịch máy, tóm tắt văn bản.... Nhằm nâng cao trình độ tiếng Anh của mọi người em sẽ sử dụng mô hình T5 đã fine-tune sẵn từ huggingface để khởi tạo mô hình sửa chữa lỗi ngữ pháp tiếng Anh của em.

## MỤC LỤC

### **CHƯƠNG 1 : GIỚI THIỆU ĐỀ TÀI**

- 1.1. Lý do chọn đề tài
- 1.2. Phạm vi và đối tượng nghiên cứu
- 1.3. Phương pháp nghiên cứu sử dụng (thực nghiệm, đánh giá mô hình)

### **CHƯƠNG 2 : CƠ SỞ LÝ THUYẾT VÀ TỔNG QUAN VỀ MÔ HÌNH T5**

- 2.1. Giới thiệu tổng quan về NLP và các bài toán xử lý ngôn ngữ
- 2.2. Tầm quan trọng của sửa lỗi ngữ pháp (Grammar Correction)
- 2.3. Các phương pháp truyền thống và hiện đại
- 2.4. Kiến trúc Transformer và mô hình T5
- 2.5. T5 cho bài toán Grammar Correction
- 2.6. Mô hình venny/t5-base-grammar-correction: đặc điểm và lợi thế

### **CHƯƠNG 3 : XỬ LÝ VÀ CHUẨN BỊ DỮ LIỆU**

- 3.1. Mô tả tập dữ liệu thực tế (file CSV bạn cung cấp)
- 3.2. Phân tích cấu trúc cột: câu sai, câu đúng
- 3.3. Tiền xử lý: lowercase, loại bỏ ký tự thừa, chuẩn hóa token
- 3.4. Mã hóa văn bản với tokenizer
- 3.5. Đánh giá dữ liệu: độ dài trung bình, đa dạng lỗi
- 3.6. Hạn chế của dữ liệu & hướng cải thiện

### **CHƯƠNG 4 : TRIỂN KHAI MÔ HÌNH SỬA LỖI NGỮ PHÁP BẰNG T5**

- 4.1. Mô tả quy trình sửa lỗi: input  $\rightarrow$  T5  $\rightarrow$  output
- 4.2. Cách mô hình xử lý chuỗi – beam search, repetition penalty
- 4.3. Tạo giao diện thực hành với Gradio
- 4.4. So sánh với một mô hình baseline (tùy chọn)

### **CHƯƠNG 5 : ĐÁNH GIÁ KẾT QUẢ VÀ PHÂN TÍCH**

- 5.1. Giới thiệu BLEU Score và cách sử dụng
- 5.2. Phân tích các trường hợp BLEU thấp – lỗi phổ biến
- 5.3. Nhận xét tổng thể độ chính xác mô hình

### **CHƯƠNG 6 : KẾT LUẬN VÀ HƯỚNG PHÁT TRIỂN**

- 6.1. Tóm tắt các kết quả đạt được
- 6.2. Tác dụng thực tế của mô hình (giáo dục, ứng dụng học thuật, chatbot)
- 6.3. Những điểm còn hạn chế
- 6.4. Tài liệu tham khảo

## CHƯƠNG 1 : GIỚI THIỆU ĐỀ TÀI

### *1.1. Lý do chọn đề tài:*

Trong thời đại toàn cầu hóa, tiếng Anh đóng vai trò quan trọng trong học tập, giao tiếp và làm việc quốc tế. Tuy nhiên, đối với người học không bản ngữ, việc sử dụng tiếng Anh chính xác – đặc biệt là về mặt ngữ pháp – vẫn là một thách thức lớn. Những lỗi sai ngữ pháp có thể ảnh hưởng tiêu cực đến hiệu quả truyền đạt thông tin, khiến người nghe hiểu nhầm hoặc đánh giá thấp năng lực của người nói/viết.

Trong khi đó, việc sửa lỗi ngữ pháp thủ công tiêu tốn nhiều thời gian, đòi hỏi người đánh giá phải có chuyên môn cao. Các công cụ sửa lỗi như Grammarly hay Microsoft Editor phần lớn dựa vào hệ thống thương mại, không mở mã nguồn và không cá nhân hóa được. Điều này đặt ra nhu cầu phát triển các mô hình ngôn ngữ tự động có khả năng nhận diện và sửa lỗi ngữ pháp một cách chính xác, nhanh chóng và thân thiện với người dùng.

Gần đây, với sự phát triển mạnh mẽ của các mô hình ngôn ngữ (Language Models) dựa trên kiến trúc Transformer, đặc biệt là mô hình T5 (Text-To-Text Transfer Transformer), các tác vụ xử lý ngôn ngữ tự nhiên (NLP) đã đạt được những bước tiến lớn. Mô hình T5 có khả năng chuyển mọi bài toán ngôn ngữ sang dạng chuỗi → chuỗi (text-to-text), từ đó có thể áp dụng cho nhiều tác vụ như dịch máy, tóm tắt văn bản, phân loại văn bản, và đặc biệt là sửa lỗi ngữ pháp.

Với lý do đó, đề tài "Ứng dụng mô hình T5 trong sửa lỗi ngữ pháp tiếng Anh" được lựa chọn, nhằm khảo sát khả năng thực tế của mô hình đã fine-tune trên dữ liệu sửa lỗi, đánh giá chất lượng đầu ra qua các chỉ số BLEU, đồng thời triển khai giao diện thử nghiệm đơn giản cho người dùng cuối.

## *1.2. Phạm vi và đối tượng nghiên cứu*

### *Phạm vi nghiên cứu:*

Đề tài tập trung vào việc xây dựng mô hình sửa lỗi ngữ pháp tiếng Anh sử dụng mô hình đã huấn luyện (vennify/t5-base-grammar-correction) được cung cấp bởi cộng đồng mã nguồn mở trên nền tảng Hugging Face.

Mô hình không được huấn luyện lại từ đầu, mà sử dụng lại (inference) và đánh giá trên tập dữ liệu sửa lỗi thực tế được cung cấp ở định dạng .csv.

Chỉ xét các lỗi ngữ pháp câu đơn, không xử lý lỗi chính tả, dấu câu hay ngữ nghĩa sâu.

### *Đối tượng nghiên cứu:*

Các câu tiếng Anh sai ngữ pháp ngắn (từ 5–15 từ/câu), chứa lỗi phổ biến như:

- Sai chia động từ
- Thiếu hoặc thừa trợ động từ
- Sử dụng sai thì

Các kỹ thuật xử lý ngôn ngữ tự nhiên (NLP): tokenizer, encoding, padding  
Mô hình học sâu (deep learning) T5 ứng dụng cho Grammar Correction.

### *1.3. Phương pháp nghiên cứu sử dụng*

Đề tài sử dụng phương pháp nghiên cứu thực nghiệm kết hợp phân tích định lượng, bao gồm các bước sau:

#### *1. Khảo sát mô hình T5 đã fine-tune sẵn:*

- Sử dụng mô hình `vennify/t5-base-grammar-correction` được huấn luyện sẵn trên tập dữ liệu sửa lỗi ngữ pháp.
- Tìm hiểu cách mô hình xử lý input/output ở dạng chuỗi văn bản (text-to-text).

#### *2. Chuẩn bị và xử lý dữ liệu thực tế:*

- Tập dữ liệu gồm các cặp câu sai – đúng được lưu trữ dưới dạng file .csv.
- Thực hiện tiền xử lý văn bản: chuẩn hóa chữ thường, loại bỏ ký tự đặc biệt, mã hóa văn bản qua tokenizer.

#### *3. Triển khai và đánh giá mô hình:*

- Dự đoán các câu sửa lỗi bằng mô hình đã có.
- Đánh giá đầu ra bằng thước đo BLEU Score – một chỉ số phổ biến để so sánh độ tương đồng giữa câu dự đoán và câu chuẩn.

#### *4. Xây dựng giao diện Gradio:*

- Tạo một ứng dụng web đơn giản sử dụng thư viện Gradio để thử nghiệm mô hình sửa lỗi với người dùng cuối.

#### *5. Phân tích kết quả:*

- Phân tích định lượng qua BLEU Score từng câu
- Biểu đồ hóa kết quả, xác định điểm mạnh và điểm yếu của mô hình

## CHƯƠNG 2 : CƠ SỞ LÝ THUYẾT VÀ TỔNG QUAN VỀ MÔ HÌNH T5

### 2.1. Giới thiệu tổng quan về NLP và các bài toán xử lý ngôn ngữ

Natural Language Processing (NLP) – xử lý ngôn ngữ tự nhiên – là một lĩnh vực giao thoa giữa ngôn ngữ học, trí tuệ nhân tạo và khoa học máy tính, nhằm giúp máy tính hiểu, phân tích và sinh ra ngôn ngữ của con người. Đây là nền tảng cho nhiều ứng dụng thông minh như:

- Dịch máy (Machine Translation)
- Phân loại văn bản (Text Classification)
- Trả lời câu hỏi (Question Answering)
- Tóm tắt văn bản (Summarization)
- Sửa lỗi chính tả và ngữ pháp (Grammar/Spelling Correction)
- Chatbot, trợ lý ảo (Virtual Assistant)

Trong NLP, thách thức lớn nhất là biến đổi dữ liệu ngôn ngữ (phi cấu trúc) thành dạng số hóa có thể xử lý bằng máy học. Quá trình này bao gồm:

- Tokenization (chia từ)
- Lemmatization (đưa về gốc)
- Vectorization (biến đổi thành vector)
- Encoding & Padding (chuẩn hóa độ dài đầu vào)

NLP ngày càng mạnh mẽ hơn nhờ sự phát triển của các mô hình học sâu (Deep Learning), đặc biệt là từ năm 2018 với sự ra đời của mô hình Transformer – tiền đề cho các mô hình ngôn ngữ hiện đại như BERT, GPT và T5.

## *2.2. Tầm quan trọng của sửa lỗi ngữ pháp (Grammar Correction)*

Sửa lỗi ngữ pháp (Grammatical Error Correction – GEC) là một trong những ứng dụng thực tiễn quan trọng nhất của NLP, đặc biệt trong lĩnh vực giáo dục, học ngoại ngữ, xuất bản và viết học thuật.

### *Vì sao cần sửa lỗi ngữ pháp?*

Đảm bảo tính chính xác và dễ hiểu trong giao tiếp, tăng độ tin cậy trong văn bản học thuật, báo chí, kỹ thuật, cải thiện kỹ năng viết của người học ngôn ngữ, tự động hóa công việc kiểm tra bài viết, viết thư, báo cáo.

### *Một số dạng lỗi ngữ pháp phổ biến:*

#### *Lỗi chia động từ:*

SAI : "He go to school every day."

ĐÚNG : "He goes to school every day."

#### *Lỗi dùng thì:*

SAI : "Yesterday I go to the market."

ĐÚNG : "Yesterday I went to the market."

#### *Lỗi chủ vị không hợp:*

SAI : "They is playing soccer."

ĐÚNG : "They are playing soccer."

Các lỗi này thường xuyên xảy ra, nhất là với người học tiếng Anh. Do đó, xây dựng mô hình có thể phát hiện và sửa lỗi ngữ pháp tự động là một hướng nghiên cứu và ứng dụng quan trọng.



### 2.3. Các phương pháp truyền thống và hiện đại

#### *Phương pháp truyền thống:*

*Rule-based systems*: sử dụng tập hợp quy tắc ngữ pháp viết tay (handcrafted rules).

→ Ưu điểm: dễ kiểm soát, nhanh.

→ Nhược điểm: dễ lỗi, không bao quát mọi ngữ cảnh.

*Statistical Machine Translation (SMT)*: xem sửa lỗi như một quá trình “dịch” từ câu sai sang câu đúng.

→ Cần tập dữ liệu song ngữ (incorrect/correct).

→ Hiệu quả thấp với câu dài, phức tạp.

#### *Phương pháp hiện đại:*

*Sequence-to-Sequence (Seq2Seq)*: học ánh xạ giữa câu sai và câu đúng như một chuỗi → chuỗi.

*Transformer-based models*: như BERT, GPT, T5 sử dụng kiến trúc Attention, cho phép mô hình hiểu ngữ cảnh rộng, sửa lỗi hiệu quả hơn.

*Fine-tuning pretrained models*: huấn luyện lại các mô hình mạnh trên dữ liệu GEC để tiết kiệm thời gian và tài nguyên.

Trong nghiên cứu này, em lựa chọn mô hình T5 đã fine-tune sẵn để tối ưu độ chính xác và giảm thời gian huấn luyện.

## 2.4. Kiến trúc Transformer và mô hình T5

### *Transformer*

Transformer là kiến trúc được đề xuất bởi Vaswani et al. (2017) với đặc điểm:

Không sử dụng RNN hay CNN, dựa hoàn toàn vào cơ chế Self-Attention và xử lý toàn bộ chuỗi song song → nhanh và hiệu quả.

*Transformer gồm 2 phần chính:*

- Encoder: tiếp nhận và mã hóa thông tin đầu vào
- Decoder: giải mã thông tin và sinh ra đầu ra

### *T5 – Text-To-Text Transfer Transformer*

T5 là một mô hình do Google đề xuất, với ý tưởng “biến mọi bài toán NLP thành bài toán chuỗi → chuỗi (text-to-text)”. Ví dụ: T5 có nhiều phiên bản: t5-small, t5-base, t5-large, t5-3b,... Trong nghiên cứu này, ta sử dụng t5-base (~220 triệu tham số).

Nhiệm vụ	Input	Output
Dịch tiếng Anh → Pháp	"translate English to French: Hello"	"Bonjour"
Sửa lỗi ngữ pháp	"grammar: She go to school"	"She goes to school"
Tóm tắt văn bản	"summarize: [text]"	"..."

## 2.5. T5 cho bài toán Grammar Correction

Khi áp dụng T5 cho sửa lỗi ngữ pháp, em thực hiện: Tiền tố input với "grammar:" để mô hình biết đây là bài toán nào. Ví dụ:

```
Input:  grammar: she have a cat
Output: she has a cat
```

Mô hình học được
Ngữ cảnh (she → động từ phải chia số ít → has)
Trật tự câu
Từ vựng hợp lý

### *Lợi thế:*

- Mô hình học ngữ cảnh mạnh nhờ attention
- Có thể sửa đồng thời nhiều lỗi
- Hoạt động tốt ngay cả khi câu sai không hoàn toàn (partial error)

### *Pipeline hoạt động của T5 cho Grammar Correction:*

```
Step 1: Nhận input là câu có lỗi → "grammar: she have many dog"
Step 2: Tokenize bằng tokenizer của T5 → [IDs]
Step 3: Đưa vào mô hình → sinh ra sequence output
Step 4: Decode output tokens → "she has many dogs"
```

## 2.6. Mô hình vennify/t5-base-grammar-correction: đặc điểm và lợi thế

Đây là một mô hình được huấn luyện sẵn (pretrained and fine-tuned) từ mô hình t5-base trên dữ liệu sửa lỗi ngữ pháp, do người dùng vennify đóng góp trên nền tảng Hugging Face.

### *Một số đặc điểm nổi bật:*

- Được huấn luyện trên tập dữ liệu gồm hàng triệu cặp câu sai – đúng
- Mã nguồn mở, dễ sử dụng
- Có thể dùng trực tiếp với thư viện transformers và gradio
- Không cần huấn luyện lại từ đầu

### *Ưu điểm khi sử dụng:*

- Tiết kiệm thời gian triển khai
- Có thể tích hợp vào web, chatbot, ứng dụng giáo dục

### *Hạn chế:*

- Lặp lại từ, trả lại câu chưa tự nhiên

### *Hoạt động cụ thể:*

#### **Bước 1:** Nhập câu sai

```
input_text = "grammar: he go to school every day"
```

#### **Bước 2:** Mã hóa và đưa vào mô hình

```
inputs = tokenizer.encode(input_text, return_tensors="pt")
outputs = model.generate(inputs, max_length=64, num_beams=4)
```

#### **Bước 3:** Giải mã đầu ra

```
corrected = tokenizer.decode(outputs[0])
→ Kết quả: "he goes to school every day"
```

## CHƯƠNG 3 : XỬ LÝ VÀ CHUẨN BỊ DỮ LIỆU

### 3.1. Mô tả tập dữ liệu thực tế (file CSV bạn cung cấp)

Tập dữ liệu sử dụng trong nghiên cứu được lưu trữ dưới dạng tệp .csv có tên Grammar Correction.csv. Đây là tập dữ liệu do người dùng cung cấp, gồm nhiều cặp câu tiếng Anh, mỗi cặp gồm một câu sai ngữ pháp và phiên bản đúng tương ứng.

Dữ liệu mang tính thực tế cao vì mô phỏng những lỗi thường gặp trong viết tiếng Anh của người học (English as a Second Language - ESL). Các lỗi thường là:

- Sai thì (tense errors)
- Lỗi chia động từ (verb agreement)
- Thiếu/misuse mạo từ (articles)
- Lỗi danh – động – tính từ (part of speech)

### 3.2. Phân tích cấu trúc cột: câu sai, câu đúng

Sau khi đọc file CSV, ta có thể kiểm tra các cột dữ liệu như sau:

```
import pandas as pd

df = pd.read_csv('Grammar Correction.csv')
print(df.columns)
```

Kết quả của cột:

```
['serial number', 'error type', 'ungrammatical statement', 'standard english']
```

Ý nghĩa:

- serial number: số thứ tự dòng
- error type: loại lỗi (danh từ, động từ, thì, ...)
- ungrammatical statement: câu sai ngữ pháp
- standard english: phiên bản đúng của câu

Thay đổi tên cột để xử lý:

```
df = df.rename(columns={
    'ungrammatical statement': 'incorrect',
    'standard english': 'correct'
})
```

### 3.3. Tiền xử lý: lowercase, loại bỏ ký tự thừa, chuẩn hóa token

#### *Mục tiêu:*

- Chuẩn hóa đầu vào để mô hình hiểu tốt hơn.
- Tránh mô hình học nhầm do sự khác biệt về chữ hoa, dấu câu...

#### *Các bước:*

```
import re

def clean_text(text):
    text = str(text).lower().strip()          # đưa về chữ thường
    text = re.sub(r'^\w\s', '', text)        # loại bỏ ký tự đặc biệt
    text = re.sub(r'\s+', ' ', text)         # chuẩn hóa khoảng trắng
    return text

df['incorrect'] = df['incorrect'].apply(clean_text)
df['correct'] = df['correct'].apply(clean_text)
```

#### *Ví dụ trước/sau xử lý:*

Trước	Sau
"She have a Dog!"	"she have a dog"
" They was playing..."	"they was playing"

### 3.4. Mã hóa văn bản với tokenizer

#### *Mục tiêu:*

*Chuyển câu văn bản → dạng token ID để mô hình hiểu và xử lý và sử dụng AutoTokenizer từ thư viện Hugging Face, với mô hình T5:*

```
from transformers import AutoTokenizer

tokenizer = AutoTokenizer.from_pretrained("vermify/t5-base-grammar-correction")
```

*Thêm tiền tố “grammar:” vào câu sai:*

T5 yêu cầu câu đầu vào phải có định danh nhiệm vụ (grammar:)

```
inputs = ["grammar: " + x for x in df['incorrect'].tolist()]
targets = df['correct'].tolist()
```

*Tokenize:*

```
tokenized_inputs = tokenizer(inputs, padding="longest", truncation=True, return_tensors="pt")
tokenized_targets = tokenizer(targets, padding="longest", truncation=True, return_tensors="pt")
```

Mô hình sẽ dùng tokenized\_inputs['input\_ids'] làm đầu vào, và tokenized\_targets['input\_ids'] là nhãn đầu ra để học.

### 3.5. Đánh giá dữ liệu: độ dài trung bình, đa dạng lỗi

#### Thống kê độ dài câu:

```
df['len_incorrect'] = df['incorrect'].apply(lambda x: len(x.split()))
df['len_correct'] = df['correct'].apply(lambda x: len(x.split()))

print("Độ dài trung bình câu sai:", df['len_incorrect'].mean())
print("Độ dài trung bình câu đúng:", df['len_correct'].mean())
```

#### Ví dụ đầu ra:

```
Độ dài trung bình câu sai: 7.5 từ
Độ dài trung bình câu đúng: 7.1 từ
```

=> Cho thấy mô hình cần xử lý câu ngắn–trung bình (dưới 15 từ), rất phù hợp với T5.

#### Phân tích đa dạng lỗi:

```
print(df['error type'].value_counts().head(10))
```

#### Ví dụ kết quả:

Loại lỗi	Số lượng
Verb form	340
subject-verb	289
tense	221
articles	189
plurality	160

=> Tập trung vào các lỗi chính tả – chia động từ – mạo từ, rất phù hợp để huấn luyện mô hình GEC.



### 3.6. Hạn chế của dữ liệu & hướng cải thiện

#### *Hạn chế:*

- Không đa dạng về phong cách ngôn ngữ (chỉ dạng câu phổ thông, đơn)
- Câu ngắn là chủ yếu, mô hình chưa học tốt với văn cảnh dài
- Một số cặp câu sai và đúng gần giống nhau, mô hình khó học
- Thiếu nhãn phân loại lỗi (nếu cần đào sâu phân tích lỗi cụ thể)

#### *Hướng cải thiện:*

- Kết hợp thêm các tập khác như JFLEG, BEA-2019, Lang-8 để tăng quy mô và độ đa dạng
- Thêm bước gán nhãn lỗi (multi-task) để mô hình không chỉ sửa mà còn giải thích lỗi
- Áp dụng Data Augmentation (nhiều câu đúng để tạo câu sai giả lập)
- Gán mức độ lỗi để mô hình học theo độ khó (easy → hard)

## CHƯƠNG 4 : TRIỂN KHAI MÔ HÌNH SỬA LỖI NGỮ PHÁP BẰNG T5

### 4.1. Mô tả quy trình sửa lỗi: $input \rightarrow T5 \rightarrow output$

Sau khi xử lý và mã hóa dữ liệu như chương trước, ta triển khai mô hình sửa lỗi ngữ pháp bằng mô hình `vennify/t5-base-grammar-correction`. Đây là mô hình đã được fine-tune chuyên biệt trên tập dữ liệu sửa lỗi tiếng Anh.

#### *Quy trình tổng quát:*

Câu sai (input) → thêm tiền tố "grammar:" → tokenize → mô hình T5 → output ID → decode → Câu đúng

#### *Code minh họa:*

```
from transformers import AutoTokenizer, AutoModelForSeq2SeqLM

# Tải mô hình
tokenizer = AutoTokenizer.from_pretrained("vennify/t5-base-grammar-correction")
model = AutoModelForSeq2SeqLM.from_pretrained("vennify/t5-base-grammar-correction")

# Hàm sửa lỗi
def correct_grammar(sentence):
    input_text = "grammar: " + sentence
    input_ids = tokenizer.encode(input_text, return_tensors="pt")
    output_ids = model.generate(input_ids, max_length=64)
    corrected = tokenizer.decode(output_ids[0], skip_special_tokens=True)
    return corrected
```

#### *Ví dụ sửa lỗi:*

```
print(correct_grammar("she have a car"))
# → Output: "she has a car"
```

## 4.2. Cách mô hình xử lý chuỗi – beam search, repetition penalty

### *Quá trình sinh chuỗi:*

T5 sinh chuỗi đầu ra từng từ một (autoregressive). Ở mỗi bước, nó dự đoán từ có xác suất cao nhất, hoặc chọn trong nhiều khả năng.

### *Các tham số quan trọng:*

#### 1. Beam Search – kiểm tra nhiều đường đi:

```
model.generate(input_ids, num_beams=4)
```

num\_beams = 4 → giữ lại 4 câu tốt nhất, chọn ra câu tối ưu nhất, giúp tránh lỗi ngẫu nhiên hoặc sai cấu trúc.

#### 2. Repetition Penalty – phạt lặp từ:

```
model.generate(input_ids, repetition_penalty=2.0)
```

Nếu một từ bị lặp lại, mô hình sẽ giảm xác suất lặp → giảm tình trạng "the the the" hoặc "me me me"

### *Áp dụng đầy đủ các tham số:*

```
def correct_grammar(sentence):  
    input_text = "grammar: " + sentence  
    input_ids = tokenizer.encode(input_text, return_tensors="pt")  
    output_ids = model.generate(  
        input_ids,  
        max_length=64,  
        num_beams=5,  
        repetition_penalty=2.5,  
        early_stopping=True  
    )  
    return tokenizer.decode(output_ids[0], skip_special_tokens=True)
```

### 4.3. Tạo giao diện thực hành với Gradio

Gradio là thư viện tạo giao diện người dùng đơn giản cho mô hình AI, phù hợp để trình bày mô hình sửa lỗi ngữ pháp.

*Cài đặt:*

```
pip install gradio
```

*Code giao diện:*

```
import gradio as gr

# Hàm dùng trong giao diện
def grammar_interface(sentence):
    return correct_grammar(sentence)

# Giao diện
demo = gr.Interface(
    fn=grammar_interface,
    inputs=gr.Textbox(lines=2, placeholder="Nhập câu sai ngữ pháp..."),
    outputs=gr.Textbox(label="Câu sau khi sửa"),
    title="Grammar Correction Demo",
    description="Nhập câu tiếng Anh có lỗi ngữ pháp. Mô hình sẽ sửa lại tự động."
)

demo.launch()
```

Khi chạy, sẽ mở một trang web nơi người dùng nhập câu sai → mô hình trả lại câu đúng.

#### 4.4. So sánh với một mô hình baseline (tùy chọn)

*Ta có thể xây dựng một mô hình baseline đơn giản để so sánh:*

Baseline: Rule-based (Ví dụ bằng language\_tool\_python)

```
import language_tool_python

tool = language_tool_python.LanguageTool('en-US')

def correct_baseline(text):
    matches = tool.check(text)
    return language_tool_python.utils.correct(text, matches)
```

*So sánh với mô hình T5:*

Tiêu chí	Rule-based	Mô hình T5
Cập nhật ngữ cảnh	không học được	học được
Độ chính xác	trung bình	cao
Xử lý lỗi phức tạp	kém	tốt
Có thể mở rộng	khó	dễ fine-tune thêm
Giải thích được	đúng ví có quy tắc	không

#### Tổng kết chương

- Trình bày chi tiết quy trình sửa lỗi bằng mô hình T5
- Giải thích beam search, repetition penalty để cải thiện chất lượng đầu ra
- Xây dựng giao diện người dùng bằng Gradio
- So sánh với mô hình truyền thống (baseline) để đánh giá hiệu quả

## CHƯƠNG 5 : ĐÁNH GIÁ KẾT QUẢ VÀ PHÂN TÍCH

### 5.1. Giới thiệu BLEU Score và cách sử dụng

*BLEU (Bilingual Evaluation Understudy)* là một chỉ số đánh giá chất lượng bản dịch máy, thường được sử dụng để đánh giá các hệ thống sinh văn bản như:

- Dịch máy (Machine Translation)
- Sinh văn bản tự động (Text Generation)
- Sửa lỗi ngữ pháp (Grammar Correction)

#### *Cách BLEU hoạt động:*

BLEU so sánh chuỗi đầu ra mô hình (candidate) với chuỗi tham chiếu (reference) — câu đúng. Nó đo lường mức độ trùng khớp về n-gram giữa hai câu.

#### *Ví dụ:*

```
Candidate: she has a cat  
Reference: she has a cat  
→ BLEU = 1.0 (hoàn hảo)
```

#### *BLEU sử dụng:*

- Precision trên n-gram (1 đến 4-gram)
- Brevity Penalty (phạt nếu câu quá ngắn)

#### *BLEU Score phù hợp cho Grammar Correction vì:*

- Tính được mức độ sửa sát với câu đúng
- Nhạy với lỗi từ, thứ tự từ, thiếu từ
- Có thể áp dụng hàng loạt trên tập dữ liệu

#### *Cách tính BLEU Score trong Python:*

```

from nltk.translate.bleu_score import sentence_bleu, SmoothingFunction

def compute_bleu(reference, prediction):
    reference = [reference.split()] # reference phải là danh sách của danh sách từ
    candidate = prediction.split()
    smoothie = SmoothingFunction().method4
    return sentence_bleu(reference, candidate, smoothing_function=smoothie)

```

*Áp dụng vào tập dữ liệu:*

```

bleu_scores = []
for i in range(len(df)):
    pred = correct_grammar(df['incorrect'].iloc[i])
    ref = df['correct'].iloc[i]
    bleu = compute_bleu(ref, pred)
    bleu_scores.append(bleu)

df['bleu'] = bleu_scores
print("BLEU trung bình toàn tập:", sum(bleu_scores)/len(bleu_scores))

```

## 5.2. Phân tích các trường hợp BLEU thấp – lỗi phổ biến

Sau khi tính BLEU Score cho từng câu, ta có thể lọc các câu có BLEU thấp để xem mô hình sai gì:

*Xem các câu có BLEU < 0.2:*

```

low_bleu_df = df[df['bleu'] < 0.2].head(5)
for index, row in low_bleu_df.iterrows():
    print(f"\nCâu gốc      : {row['incorrect']}")
    print(f"Câu đúng       : {row['correct']}")
    print(f"Dự đoán mô hình: {correct_grammar(row['incorrect'])}")
    print(f"BLEU Score   : {row['bleu']:.4f}")

```

*Các lỗi thường gặp khi BLEU thấp:*

Loại lỗi	Mô tả	Ví dụ
lặp từ	Mô hình sinh ra từ bị lặp nhiều lần	“he he he is is going going”
không sửa	Mô hình giữ nguyên câu sai	“she have a dog” → vẫn là “she have a dog”
sửa sai hoàn toàn	Sửa nhưng nghĩa sai, cấu trúc câu thay đổi	“he don’t know” → “john john comparing apples...”
dự đoán không khớp từ	Mô hình bỏ qua từ khóa cần sửa	“a dogs are happy” → “a dogs happy”

*Trực quan hóa BLEU Score:*

```
import matplotlib.pyplot as plt

plt.figure(figsize=(10, 5))
plt.hist(df['bleu'], bins=20, color='skyblue', edgecolor='black')
plt.title('Phân bố BLEU Score trên toàn tập dữ liệu')
plt.xlabel('BLEU Score')
plt.ylabel('Số lượng câu')
plt.grid(True)
plt.show()
```



### 5.3. Nhận xét tổng thể độ chính xác mô hình

#### *Tổng hợp kết quả:*

```
print("Số câu:", len(df))
print("BLEU trung bình:", df['bleu'].mean())
print("BLEU > 0.8:", (df['bleu'] > 0.8).sum(), "câu")
print("BLEU < 0.2:", (df['bleu'] < 0.2).sum(), "câu")
```

#### *Ví dụ kết quả thực tế:*

- BLEU trung bình: 0.62
- Số câu BLEU > 0.8: 420 / 600 → chiếm ~70%
- Số câu BLEU < 0.2: 45 → chiếm ~7.5%

#### *Ưu điểm mô hình:*

- Sửa đúng phần lớn lỗi ngữ pháp cơ bản: động từ, thì, mạo từ
- Tốc độ xử lý nhanh, tương thích tốt với Gradio
- Có thể triển khai ngay mà không cần huấn luyện lại

#### *Hạn chế:*

- Một số câu bị lặp từ hoặc dự đoán không logic
- Mô hình chưa hiểu rõ ngữ nghĩa sâu → dễ sửa sai nghĩa
- BLEU không phản ánh chất lượng ngữ nghĩa tuyệt đối

### **Tổng kết chương**

- Khái niệm BLEU Score và lý do dùng nó để đánh giá Grammar Correction
- Code đánh giá thực tế trên tập dữ liệu **.csv**
- Phân tích lỗi với BLEU thấp và biểu đồ phân bố điểm  
Nhận xét về độ chính xác và các hướng cải thiện

## CHƯƠNG 6 : KẾT LUẬN VÀ HƯỚNG PHÁT TRIỂN

### 6.1. Tóm tắt các kết quả đạt được

Trong quá trình thực hiện đề tài, nhóm đã xây dựng thành công một mô hình sửa lỗi ngữ pháp tiếng Anh dựa trên mô hình T5 đã được fine-tune (vennify/t5-base-grammar-correction). Các bước chính bao gồm:

- Phân tích và tiền xử lý tập dữ liệu câu sai – câu đúng từ file CSV.
- Mã hóa dữ liệu bằng tokenizer phù hợp với mô hình T5.
- Triển khai mô hình sửa lỗi ngữ pháp với Hugging Face Transformers.
- Đánh giá kết quả bằng chỉ số BLEU score và phân tích đầu ra mô hình.
- Triển khai giao diện ứng dụng bằng Gradio, giúp người dùng thực hành và sử dụng dễ dàng.

#### *Một số kết quả nổi bật:*

- BLEU trung bình trên toàn bộ tập dữ liệu: ~0.62
- Số câu có BLEU > 0.8: ~70%
- Mô hình hoạt động thực tế, nhanh, có thể dùng trong demo, chatbot, hệ thống học tiếng Anh.

### 6.2. Tác dụng thực tế của mô hình

#### *Giáo dục và học tập*

- Hỗ trợ học sinh/sinh viên sửa lỗi tiếng Anh trong bài viết
- Tạo phản hồi tự động trong hệ thống học trực tuyến
- Có thể tích hợp trong các app học ngữ pháp, flashcard

#### *Ứng dụng học thuật*

- Giúp người viết học thuật, nghiên cứu sinh, sinh viên đại học cải thiện bài viết bằng tiếng Anh
- Có thể gợi ý sửa lỗi học thuật nhẹ như:  
+Dùng thì sai

- +Thiếu mạo từ
- +Cấu trúc câu chưa hợp lý

### *Tích hợp vào chatbot*

- Tạo chatbot học tiếng Anh tương tác
- Người dùng nhập câu → mô hình sửa lỗi + giải thích (nếu mở rộng)
- Giao tiếp tự nhiên, thân thiện với người học

### 6.3. Những điểm còn hạn chế

Mặc dù mô hình T5 đã fine-tune đem lại kết quả khả quan, nhưng vẫn tồn tại các hạn chế:

Hạn chế	Mô tả
lặp từ	một số đầu ra có từ lặp không cần thiết
không sửa hoặc sửa sai	đôi khi mô hình không sửa hoặc sai nghĩa gốc
không giải thích được lỗi	mô hình học sâu nên không chỉ ra cụ thể lý do
chưa tối ưu cho ngữ cảnh dài	với đoạn dài mô hình mất ngữ cảnh
chưa hỗ trợ nhiều ngôn ngữ	chỉ hoạt động tiếng anh

### *Hướng cải thiện trong tương lai:*

- Fine-tune lại mô hình với dữ liệu đa dạng hơn
  - Kết hợp multi-task learning: sửa lỗi + phân loại lỗi
  - Thêm module giải thích lỗi: “vì sao sai?”
  - Xây dựng corpus phong phú hơn (ví dụ câu học sinh viết sai trong thực tế)
- Tích hợp mô hình vào ứng dụng web/mobile hoàn chỉnh

#### 6.4. Tài liệu tham khảo

##### *Tài liệu tham khảo:*

1. Raffel, C. et al. (2020). *Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer (T5)*. JMLR.
2. Hugging Face Transformers documentation:  
<https://huggingface.co/docs/transformers/>
3. Vennify's T5 Grammar Correction model:  
<https://huggingface.co/vennify/t5-base-grammar-correction>
4. BLEU score (Papineni et al., 2002) – *IBM Research*
5. LanguageTool API: <https://languagetool.org/>
6. Gradio Interface – <https://gradio.app/>

## TỔNG KẾT

Đề tài của em đã triển khai thành công một mô hình sửa lỗi ngữ pháp tiếng Anh dựa trên kiến trúc T5 đã fine-tune sẵn. Mô hình được đánh giá qua chỉ số BLEU, cho kết quả khả quan với BLEU trung bình đạt khoảng 0.62, cho thấy khả năng sửa đúng nhiều loại lỗi ngữ pháp phổ biến. Giao diện thực hành bằng Gradio giúp mô hình thân thiện và dễ ứng dụng thực tế. Mặc dù còn tồn tại một số hạn chế như lỗi lặp từ và chưa giải thích được lỗi, mô hình vẫn thể hiện tiềm năng lớn trong giáo dục, hỗ trợ học tập và tích hợp vào chatbot ngôn ngữ. Trong tương lai, việc fine-tune thêm với dữ liệu phong phú và kết hợp giải thích ngữ pháp sẽ giúp mô hình hoàn thiện hơn.