

4) Insurance Multilingual Policy Document Dataset

Project Observations:

- **Data Diversity:** The dataset includes multilingual insurance policies, providing a rich resource for cross-language NLP tasks.
- **Text Complexity:** Policy texts vary in length and structure, requiring robust models for accurate translation and summarization.
- **Model Performance:** The fine-tuned mBART model shows strong potential for generating concise, high-quality summaries across languages.
- **Challenges:** Handling language-specific nuances and maintaining consistency in summarization quality remain key areas for improvement.
- **Future Scope:** Further tuning, error analysis, and exploring other transformer models (like mT5) can enhance performance and deployment readiness.

Index

- 1 - EDA - [4 - 9 Slide]
- 2 - Preprocessing - [11 - 26 Slide]
- 3 - Model Training - [27 - 30 Slide]
- 4 - mBART - [31 - 45 Slide]
- 5 - mT5 - [46 - 53 Slide]
- 6 - Save Model - [54 - 55 Slide]
- 7 - Load Model & Prediction - [56 - 71 Slide]
- 8 - Prediction - [72 - 73 Slide]
- 9 - Streamlit ChatBot - [74 - 75 Slide]
- 10 - Conclusion - [76 - 77 Slide]

Using Real Time Dataset:

- 1]Newspaper Article Dataset From Kaggle [77 - 84 Slide]
- 2] Text Preprocessing [85 - 88 Slide]
- 3] mBART / mT5 Auto tokenizer [91 - 93 Slide]
- 4] Fine Tuning for mBART [94 - 98 Slide]
- 5] Fine Tuning for mT5 [99 - 105 Slide]
- 6] Try Again FIne Tuning for mBART- Large [106 - 107 Slide]
- 7] Fine Tuning for BART-small [108 - 111 Slide]
- 8] Prediction [112 Slide]
- 9] Streamlit [113 - 114 Slide]
- 10] Machine Translation [115 - 128 Slide]
- 11] Conclusion [129 Slide]

EDA

df4 - Insurance Multilingual Policy Document Dataset



Description: Collection of insurance policies in multiple languages for NLP-based translation and summarization.

```
[1] 1 import pandas as pd  
  
[2] 1 df4=pd.read_csv('content/df4-insurance_multilingual_policy_document_dataset.csv')
```

▶ 1 df4.head()

	Unnamed: 0	Policy_ID	Policy_Text	Summarized_Text
0	0	bdd640fb-0667-4ad1-9c80-317fa3b1799d	Development say quality throughout beautiful. ...	Detail food shoulder argue start source husban...
1	1	1c11f735-dc71-4d96-8c0f-d195c17af08a	Civil quite others his other life edge network...	She campaign little near enter their instituti...
2	2	5496f63c-dc11-40c1-880a-adfbe7c99b26	Draw protect Democrat car very number line. Sp...	Total clearly able hospital unit size expect r...
3	3	14822f53-8201-4c62-b5f5-9b220e8fa8e0	Organization push dog build. East organization...	Data plant enough major town suffer begin inte...
4	4	5e84f058-d5a8-44eb-8939-23de8babce3b	Challenge camera final together someone team t...	Image street fight decision size parent focus ...

```
1 df4.nunique()
```

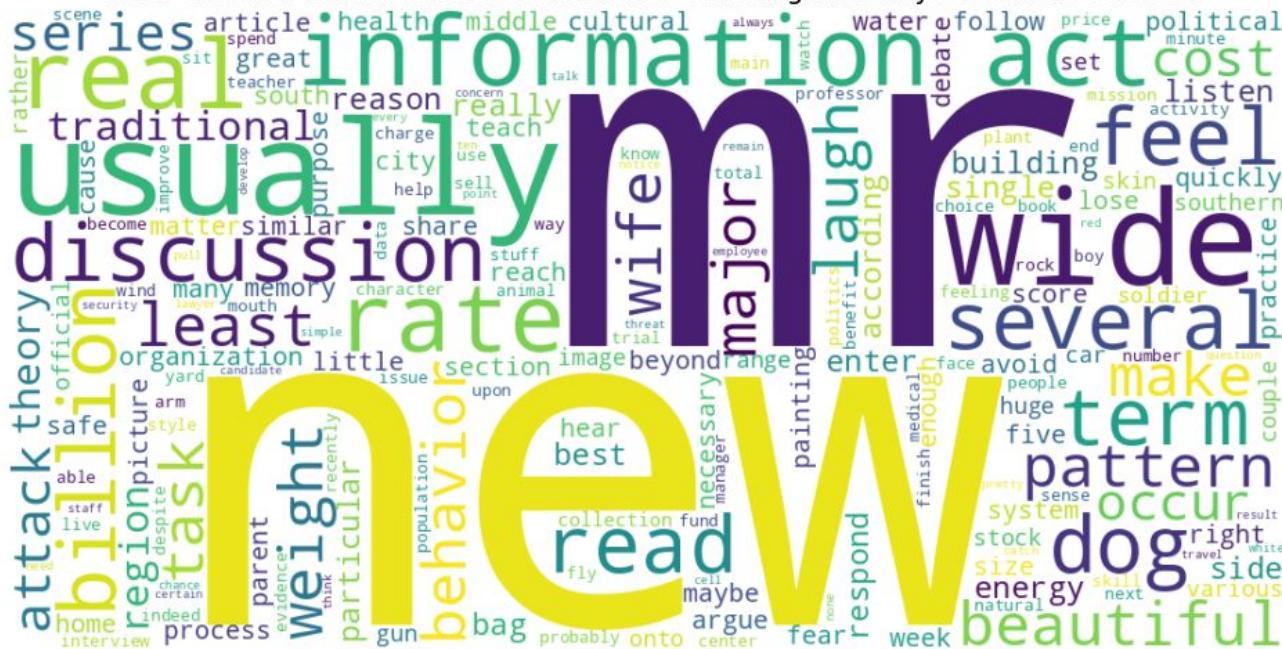
0

Policy_ID	1000
Policy_Text_EN	10
Policy_Text_FR	10
Policy_Text_ES	10
Summarized_Text	10

Create wordcloud for TEXT column for better Understanding

▼ 1) Policy_Text

Word Cloud of Textual Data from Insurance Multilingual Policy Document Dataset



▼ 2) Summarized_Text

Word Cloud of Textual Data from Insurance Multilingual Policy Document Dataset



NEXT STEP: TEXT Preprocessing

- 1) Tokenization
- 2) Stop Word Removal
- 3) Lowercasing
- 4) Remove Punctuation
- 5) Lemmatization
- 6) Vectorization

Feature Extraction:

- 1) Bag of Words
- 2) TF - IDF (TF - Term Frequency, IDF - Inverse Document Frequency)
- 3) Word Embedding
- 4) N-Grams
- 5) POS (Part of Speech)

Preprocessing

```
1 df4.head()
```

	Policy_ID	Policy_Text_EN	Policy_Text_FR	Policy_Text_ES	Summarized_Text
0	bdd640fb-0667-4ad1-9c80-317fa3b1799d	Business insurance protects against liability ...	L'assurance habitation protège contre les domm...	La póliza de salud asegura la cobertura de gas...	Protects home from natural disasters.
1	23b8c1e9-3924-46de-beb1-3b9046685257	Business insurance protects against liability ...	L'assurance pour animaux couvre les frais vété...	Esta póliza cubre los daños a su vehículo caus...	Covers vehicle accident damage.
2	bd9c66b3-ad3c-4d6d-9a3d-1fa7bc8960a9	Vehicle insurance covers theft and vandalism.	L'assurance médicale couvre les frais d'hospit...	El seguro de vida proporciona seguridad financ...	Covers hospitalization costs.
3	972a8469-1641-4f82-8b9d-2434e465e150	Travel insurance covers trip cancellations and...	L'assurance entreprise protège contre les récl...	El seguro médico cubre los gastos de hospitali...	Includes fire damage protection.
4	17fc695a-07a0-4a6e-8822-e8f36c031199	Vehicle insurance covers theft and vandalism.	L'assurance automobile couvre le vol et le van...	El seguro médico cubre los gastos de hospitali...	Covers vehicle accident damage.

✓ 1) Text Normalization

- 1) Convert text to lowercase,
- 2) remove special characters,
- 3) expand contractions. - Change the Size of the Word

i) Convert text to lowercase

✓ Policy_Text & Summarized_Text

```
[ ] 1 df4['Policy_Text_EN'] = df4['Policy_Text_EN'].str.lower()  
2 df4['Policy_Text_FR'] = df4['Policy_Text_FR'].str.lower()  
3 df4['Policy_Text_ES'] = df4['Policy_Text_ES'].str.lower()  
4 df4['Summarized_Text'] = df4['Summarized_Text'].str.lower()
```

```
[ ] 1 df4.head()
```

	Policy_ID	Policy_Text_EN	Policy_Text_FR	Policy_Text_ES	Summarized_Text
0	bdd640fb-0667-4ad1-9c80-317fa3b1799d	business insurance protects against liability ...	l'assurance habitation protège contre les domm...	la póliza de salud asegura la cobertura de gas...	protects home from natural disasters.
1	23b8c1e9-3924-46de-beb1-3b9046685257	business insurance protects against liability ...	l'assurance pour animaux couvre les frais vété...	esta póliza cubre los daños a su vehículo caus...	covers vehicle accident damage.
2	bd9c66b3-ad3c-4d6d-9a3d-1fa7bc8960a9	vehicle insurance covers theft and vandalism.	l'assurance médicale couvre les frais d'hospit...	el seguro de vida proporciona seguridad financ...	covers hospitalization costs.

▼ ii) remove special characters

```
[ ] 1 import re  
2 df4['Policy_Text_EN'] = df4['Policy_Text_EN'].str.lower().apply(lambda x: re.sub(r'[^a-zA-Z0-9\s]', '', x))  
3 df4['Policy_Text_FR'] = df4['Policy_Text_FR'].str.lower().apply(lambda x: re.sub(r'[^a-zA-Z0-9\s]', '', x))  
4 df4['Policy_Text_ES'] = df4['Policy_Text_ES'].str.lower().apply(lambda x: re.sub(r'[^a-zA-Z0-9\s]', '', x))  
5 df4['Summarized_Text'] = df4['Summarized_Text'].str.lower().apply(lambda x: re.sub(r'[^a-zA-Z0-9\s]', '', x))
```

```
[ ] 1 df4.head()
```

	Policy_ID	Policy_Text_EN	Policy_Text_FR	Policy_Text_ES	Summarized_Text
0	bdd640fb-0667-4ad1-9c80-317fa3b1799d	business insurance protects against liability ...	assurance habitation protge contre les dommag...	la pliza de salud asegura la cobertura de gast...	protects home from natural disasters
1	23b8c1e9-3924-46de-beb1-3b9046685257	business insurance protects against liability ...	assurance pour animaux couvre les frais vtrin...	esta pliza cubre los daos a su vehculo causado...	covers vehicle accident damage
2	bd9c66b3-ad3c-4d6d-9a3d-1fa7bc8960a9	vehicle insurance covers theft and vandalism	assurance mdicale couvre les frais dhospitali...	el seguro de vida proporciona seguridad financ...	covers hospitalization costs
3	972a8469-1641-4f82-8b9d-	travel insurance covers trip	assurance entreprise protge contre	el seguro mdico cubre los gastos de	includes fire damage

✓ iii) expand contractions.

```
[ ] 1 !pip install contractions
```

```
[ ] 1 import re
2 import contractions
3
4 df4['Policy_Text_EN'] = df4['Policy_Text_EN'].str.lower().apply(lambda x: re.sub(r'[^a-zA-Z0-9\s]', '', contractions.fix(x)))
5 df4['Policy_Text_FR'] = df4['Policy_Text_FR'].str.lower().apply(lambda x: re.sub(r'[^a-zA-Z0-9\s]', '', contractions.fix(x)))
6 df4['Policy_Text_ES'] = df4['Policy_Text_ES'].str.lower().apply(lambda x: re.sub(r'[^a-zA-Z0-9\s]', '', contractions.fix(x)))
7 df4['Summarized_Text'] = df4['Summarized_Text'].str.lower().apply(lambda x: re.sub(r'[^a-zA-Z0-9\s]', '', contractions.fix(x)))
8
```

```
▶ 1 df4.head()
```

	Policy_ID	Policy_Text_EN	Policy_Text_FR	Policy_Text_ES	Summarized_Text
0	bdd640fb-0667-4ad1-9c80-317fa3b1799d	business insurance protects against liability ...	assurance habitation protge contre les dommag...	la pliza de salud asegura la cobertura de gast...	protects home from natural disasters
1	23b8c1e9-3924-46de-beb1-3b9046685257	business insurance protects against liability ...	assurance pour animaux couvre les frais vtrin...	esta pliza cubre los daos a su vehculo causado...	covers vehicle accident damage
2	bd9c66b3-ad3c-4d6d-9a3d-1fa7bc8960a9	vehicle insurance covers theft and vandalism	assurance mdicale couvre les frais dhospitali...	el seguro de vida proporciona seguridad financ...	covers hospitalization costs
3	972a8469-1641-4f82-8b9d-2434e465e150	travel insurance covers trip cancellations and...	assurance entreprise protge contre les reclama...	el seguro mdico cubre los gastos de hospitaliz...	includes fire damage protection

✓ expand contractions example:

```
import contractions

text = "I'm going to the store."

expanded_text = contractions.fix(text)

print(expanded_text) # Output: "I am going to the store."
```

2) Tokenization & Embeddings

▼ i) Tokenization

- a) Sentence Tokenization
- b) Word Tokenization
- c) Letter Tokenization

```
6 # Step 2: Reinstall a stable version of NLTK  
7 !pip install nltk==3.8.1 # Pin to a stable version
```

```
[ ]    1 import nltk  
      2 from nltk.tokenize import word_tokenize  
      3 nltk.download('punkt')
```

```
→ [nltk_data] Downloading package punkt to /root/nltk_data...  
[nltk_data]  Unzipping tokenizers/punkt.zip.  
True
```

```
[ ] 1 # Tokenizing text  
2 df4['Policy_Text_Tokens_EN'] = df4['Policy_Text_EN'].apply(word_tokenize)  
3 df4['Policy_Text_Tokens_FR'] = df4['Policy_Text_FR'].apply(word_tokenize)  
4 df4['Policy_Text_Tokens_ES'] = df4['Policy_Text_ES'].apply(word_tokenize)  
5 df4['Summarized_Text_Tokens'] = df4['Summarized_Text'].apply(word_tokenize)  
6
```

Policy_Text_Tokens_EN	Policy_Text_Tokens_FR	Policy_Text_Tokens_ES	Summarized_Text_Tokens
[Business, insurance, protects, against, liabi...]	[L'assurance, habitation, protège, contre, les...]	[La, póliza, de, salud, asegura, la, cobertura...]	[Protects, home, from, natural, disasters,]
[Business, insurance, protects, against, liabi...]	[L'assurance, pour, animaux, couvre, les, fria...]	[Esta, póliza, cubre, los, daños, a, su, vehíc...]	[Covers, vehicle, accident, damage,]
[Vehicle, insurance, covers, theft, and, vanda...]	[L'assurance, médicale, couvre, les, frais, d'...]	[El, seguro, de, vida, proporciona, seguridad,...]	[Covers, hospitalization, costs,]
[Travel, insurance, covers, trip, cancellation...]	[L'assurance, entreprise, protège, contre, les...]	[El, seguro, médico, cubre, los, gastos, de, h...]	[Includes, fire, damage, protection,]

ii) Embeddings

Embeddings Means convert text or image to vectorization

Convert Word to Number to give particular weight certain words

"An embedding is a way to convert words into numerical vectors, capturing their meaning and context. It assigns particular weights based on word usage, relationships, and importance — where meaningful, frequently used, or contextually important words often get more distinct representations."

 **Byte Pair Encoding (BPE):**

- **How it works:** BPE merges the most frequent pairs of characters or subwords iteratively to form new subwords.
- **Used by:** GPT models, RoBERTa, OpenAI models.
- **Strengths:**
 - Efficient and fast.
 - Works well for language modeling tasks.
 - Handles rare and out-of-vocabulary words well by breaking them down into subword units.
- **Limitations:**
 - Pure frequency-based, so it might not always capture the meaning behind word parts.

WordPiece Tokenization:

- **How it works:** WordPiece uses a similar merging approach but balances frequency and likelihood (how much a subword improves the language model's performance).
- **Used by:** BERT, DistilBERT, mBART, ALBERT.
- **Strengths:**
 - More optimized for model performance.
 - Often results in smaller vocabulary size with good coverage.
 - Better for handling morphologically rich languages.
- **Limitations:**
 - Slightly more complex than BPE.



For your work with BERT embeddings and mBART, WordPiece Tokenization is likely the better choice because it's optimized for transformer-based models. It balances efficiency and effectiveness,

✓ using WordPiece Tokenization

```
[ ] 1 from nltk.tokenize import word_tokenize  
2 from transformers import BertTokenizer
```

✓ Bidirectional Encoder Representations from Transformers (BERT)

```
[ ] 1 from transformers import BertTokenizer, BertModel  
2 import torch  
3  
4 # Load multilingual BERT  
5 tokenizer = BertTokenizer.from_pretrained('bert-base-multilingual-cased')  
6 model = BertModel.from_pretrained('bert-base-multilingual-cased')  
7  
8 # Function to get BERT embeddings  
9 def get_embeddings(text):  
10    inputs = tokenizer(text, return_tensors='pt', padding=True, truncation=True)  
11    with torch.no_grad():  
12        outputs = model(**inputs)  
13    return outputs.last_hidden_state.mean(dim=1).squeeze().numpy()  
14  
15 # Apply embeddings for each language and summary  
16 df4['Policy_Text_Embedding_EN'] = df4['Policy_Text_EN'].apply(get_embeddings)  
17 df4['Policy_Text_Embedding_FR'] = df4['Policy_Text_FR'].apply(get_embeddings)  
18 df4['Policy_Text_Embedding_ES'] = df4['Policy_Text_ES'].apply(get_embeddings)  
19 df4['Summarized_Text_Embedding'] = df4['Summarized_Text'].apply(get_embeddings)  
20  
21 # Save the updated DataFrame  
22 df4.to_csv('df4-embeddings-nlp.csv', index=False)  
23
```

```
→ /usr/local/lib/python3.11/dist-packages/huggingface_hub/utils/_auth.py:94: UserWarning:  
The secret `HF_TOKEN` does not exist in your Colab secrets.  
To authenticate with the Hugging Face Hub, create a token in your settings tab (https://huggingface.co/settings/tokens),  
You will be able to reuse this secret in all of your notebooks.  
Please note that authentication is recommended but still optional to access public models or datasets.  
    warnings.warn(  
tokenizer_config.json: 100%  49.0/49.0 [00:00<00:00, 2.08kB/s]  
vocab.txt: 100%  996k/996k [00:00<00:00, 4.24MB/s]  
tokenizer.json: 100%  1.96M/1.96M [00:00<00:00, 16.0MB/s]  
config.json: 100%  625/625 [00:00<00:00, 13.1kB/s]  
model.safetensors: 100%  714M/714M [00:08<00:00, 112MB/s]  
Embeddings generated and saved successfully!
```

What's happening here:

- We use `bert-base-multilingual-cased`, which supports multiple languages including English, French, and Spanish.
- `get_embeddings()` tokenizes the text, runs it through the BERT model, and takes the mean of the last hidden state to create a fixed-size embedding.
- We generate embeddings for each policy text language and the summarized text.
- Finally, we save everything in a new CSV.

3) Parallel Corpus Alignment

Use sentence alignment techniques to ensure source and translated texts match correctly

```
[ ]    1 !pip install deep-translator  
2
```

```
[ ] 1 from sklearn.metrics.pairwise import cosine_similarity
2 import numpy as np
3
4 # Function to align sentences based on cosine similarity
5 def align_sentences(source_texts, target_texts):
6     aligned_pairs = []
7     for src, tgt in zip(source_texts, target_texts):
8         src_emb = get_embeddings(src)
9         tgt_emb = get_embeddings(tgt)
10        similarity = cosine_similarity([src_emb], [tgt_emb])[0][0]
11        aligned_pairs.append((src, tgt, similarity))
12    return aligned_pairs
13
14 # Aligning English, French, and Spanish texts with their summaries
15 df4['Aligned_EN'] = df4.apply(lambda row: align_sentences([row['Policy_Text_EN']], [row['Summarized_Text']]), axis=1)
16 df4['Aligned_FR'] = df4.apply(lambda row: align_sentences([row['Policy_Text_FR']], [row['Summarized_Text']]), axis=1)
17 df4['Aligned_ES'] = df4.apply(lambda row: align_sentences([row['Policy_Text_ES']], [row['Summarized_Text']]), axis=1)
18
19 # Extracting alignment scores
20 df4['Alignment_Score_EN'] = df4['Aligned_EN'].apply(lambda x: x[0][2])
21 df4['Alignment_Score_FR'] = df4['Aligned_FR'].apply(lambda x: x[0][2])
22 df4['Alignment_Score_ES'] = df4['Aligned_ES'].apply(lambda x: x[0][2])
23
24 # Save the updated DataFrame
25 df4.to_csv('df4-aligned-embeddings.csv', index=False)
26
27 print("Sentence alignment completed and saved successfully!")
28
```

→ Sentence alignment completed and saved successfully!

Aligned_EN	Aligned_FR	Aligned_ES	Alignment_Score_EN	Alignment_Score_FR	Alignment_Score_ES
["Business insurance protects against liability risk"]	["L'assurance habitation protège contre les dommages matériels et financiers"]	["La póliza de salud asegura la cobertura de los gastos médicos y hospitalarios"]	0.629268	0.499744	0.429614
["Business insurance protects against liability risk"]	["L'assurance pour animaux couvre les frais vétérinaires"]	["Esta póliza cubre los daños a su vehículo causados por accidentes o robo"]	0.608174	0.320818	0.481441
["Vehicle insurance covers theft and vandalism risk"]	["L'assurance médicale couvre les frais d'hôpital et de soins de suite"]	["El seguro de vida proporciona seguridad financiera en caso de muerte o invalidez permanente"]	0.572594	0.438874	0.426255
["Travel insurance covers trip cancellations risk"]	["L'assurance entreprise protège contre les risques de annulation de voyage"]	["El seguro médico cubre los gastos de hospitales y farmacias en caso de enfermedad o accidente"]	0.595465	0.392005	0.440118
["Vehicle insurance covers theft and vandalism risk"]	["L'assurance automobile couvre le vol et les accidents"]	["El seguro médico cubre los gastos de hospitales y farmacias en caso de enfermedad o accidente"]	0.639730	0.495426	0.437731

What this does:

- Uses `get_embeddings()` from the previous step.
- Calculates **cosine similarity** between policy text and summary embeddings.
- Captures the alignment score (0 to 1) — closer to 1 means better alignment.
- Aligns **English (EN)**, **French (FR)**, and **Spanish (ES)** policy texts with their respective summaries.
- Saves everything in the new CSV.

Model Training

df4 - Model Training

Used to train machine translation models for multilingual document translation.

- Let's fine-tune the mBART and mT5 models on policy text dataset!

We'll train them for a text summarization task – using the multilingual policy text (Policy_Text_EN, Policy_Text_FR, Policy_Text_ES) as input and the AI-generated Summarized_Text as the target.

```
[ ] 1 pip install transformers datasets sentencepiece accelerate
```

2. Load Data & Preprocess

```
[ ] 1 import pandas as pd
2 from datasets import Dataset
3
4 # Load your aligned dataset
5 df = pd.read_csv('/content/df4-aligned-embeddings.csv')
6
7 # Keep only needed columns
8 df = df[['Policy_Text_EN', 'Summarized_Text']]
9
10 # Convert to Hugging Face Dataset format
11 dataset = Dataset.from_pandas(df)
12
```

python

```
dataset = Dataset.from_pandas(df)
```

Copy Edit

- Converts the Pandas DataFrame `df` into a **Hugging Face Dataset**.
- **Why?** Because Hugging Face models like `Trainer` expect datasets in this format.
- This format is more efficient when working with large text datasets and allows for seamless integration with tokenization and model training.

▶ 1 dataset

```
→ Dataset({  
    features: ['Policy_Text_EN', 'Summarized_Text'],  
    num_rows: 1000  
})
```

[] 1 df.head()

→

	Policy_Text_EN	Summarized_Text
0	Business insurance protects against liability ...	Protects home from natural disasters.
1	Business insurance protects against liability ...	Covers vehicle accident damage.
2	Vehicle insurance covers theft and vandalism.	Covers hospitalization costs.
3	Travel insurance covers trip cancellations and...	Includes fire damage protection.
4	Vehicle insurance covers theft and vandalism.	Covers vehicle accident damage.

mBART

▼ 3. Tokenizer Setup (for mBART)

```
[ ] 1 from transformers import MBartTokenizer  
2  
3 model_name = "facebook/mbart-large-50"  
4 tokenizer = MBartTokenizer.from_pretrained(model_name)  
5  
6 # Tokenize the dataset  
7 def preprocess_function(examples):  
8     model_inputs = tokenizer(examples['Policy_Text_EN'], max_length=512, truncation=True)  
9     labels = tokenizer(examples['Summarized_Text'], max_length=128, truncation=True)  
10    model_inputs['labels'] = labels['input_ids']  
11    return model_inputs  
12  
13 tokenized_dataset = dataset.map(preprocess_function, batched=True)  
14
```

→ The tokenizer class you load from this checkpoint is not the same type as the class this function is called from. It may result in unexpected tokenization.
The tokenizer class you load from this checkpoint is 'MBart50Tokenizer'.
The class this function is called from is 'MBartTokenizer'.

Map: 100%  1000/1000 [00:00<00:00, 2500.86 examples/s]

```
from transformers import MBartTokenizer
```

 Copy  Edit

- `MBartTokenizer` : This is the tokenizer specifically built for mBART — a multilingual sequence-to-sequence model from Hugging Face's `transformers` library.
- Tokenizers convert text into numerical representations that the model can understand — like turning sentences into tokens (numbers).

```
model_name = "facebook/mbart-large-50"  
tokenizer = MBartTokenizer.from_pretrained(model_name)
```

 Copy  Edit

- `model_name` : `"facebook/mbart-large-50"` is the pre-trained mBART model with support for 50 languages.
- `from_pretrained()` : Loads the pre-trained tokenizer associated with this model from the Hugging Face model hub.
- Now, the `tokenizer` knows how to convert words into tokens (numbers) and back.

```
def preprocess_function(examples):
    model_inputs = tokenizer(examples['Policy_Text_EN'], max_length=512) ⌂ Copy ⌂ Edit Tr
    labels = tokenizer(examples['Summarized_Text'], max_length=128, truncation=True)
    model_inputs['labels'] = labels['input_ids']
    return model_inputs
```

- `preprocess_function()` : Prepares the data for the model — converting text into tokenized format.
- `examples['Policy_Text_EN']` : The original (input) policy text.
- `examples['Summarized_Text']` : The target text — the summary we want the model to learn to generate.



- `tokenizer()`:
 - `max_length=512` : Caps the input length to 512 tokens to avoid too-long sequences.
 - `truncation=True` : If the text is longer than the max length, it gets truncated to fit.
- `labels['input_ids']` : These are the numerical token IDs representing the summary — set as the **labels** for training.
- Finally, it returns `model_inputs`, now containing:
 - Tokenized policy text.
 - Tokenized summary as labels.

```
tokenized_dataset = dataset.map(preprocess_function, batched=True) Copy Edit
```

- `map()` : Applies the `preprocess_function()` to every row in the `dataset` .
- `batched=True` : Processes the data in batches for efficiency, instead of one row at a time.
- `tokenized_dataset` : A new dataset where both inputs and labels are tokenized and ready for training.

So essentially, this whole process turns your text into a format the mBART model can use – input IDs and target IDs – so it can learn the mapping from policy text to summarized text

⌄ 4. Model & Training Setup (mBART)

The tokenizer class you load from this checkpoint is not the same type as the class this function is called from. It may result in unexpected behavior.
The tokenizer class you load from this checkpoint is 'MBart50Tokenizer'.
The class this function is called from is 'MBartTokenizer'.

Map: 100% [██████████] 30/30 [00:00<00:00, 390.98 examples/s]

```
/usr/local/lib/python3.11/dist-packages/transformers/training_args.py:1575: FutureWarning: `evaluation_strategy` is deprecated and will be removed in version 5.0.0 for `Trainer.__init__`. Use
  warnings.warn(
<ipython-input-1-07d80d70cb0c>:57: FutureWarning: `tokenizer` is deprecated and will be removed in version 5.0.0 for `Trainer.__init__`. Use
    trainer = Trainer(
[██████████] [1/1 00:00, Epoch 1/1]
```

Step Training Loss

```
import gc
import torch

gc.collect()
torch.cuda.empty_cache()
```

ⓘ Copy ⚙ Edit

- `gc.collect()` : Runs Python's garbage collector to clean up unused memory.
- `torch.cuda.empty_cache()` : Clears the GPU memory that's not actively being used — helpful when you've had crashes or memory leaks.

```
cols_needed = ['Policy_Text_EN', 'Summarized_Text']
df = pd.read_csv('/content/synthetic_policy_dataset.csv', usecols=cols_needed)
df = df.sample(frac=0.03).reset_index(drop=True)
```

- `usecols=cols_needed` : Only loads the necessary columns to avoid wasting memory.
- `df.sample(frac=0.03)` : Uses only **3% of the data** — a smart move when you keep hitting memory crashes.
- `reset_index(drop=True)` : Resets the index after sampling so the dataset stays clean.

```
small_dataset = Dataset.from_dict({
    'Policy_Text_EN': df['Policy_Text_EN'].tolist(),
    'Summarized_Text': df['Summarized_Text'].tolist()
})
```

- `Dataset.from_dict()` : Converts your Pandas DataFrame into a Hugging Face `Dataset` format, which is optimized for model training.
- `tolist()` : Converts each column to a list so it's compatible with `Dataset`.

```
model_name = 'facebook/mbart-large-50'  
tokenizer = MBartTokenizer.from_pretrained(model_name)  
model = MBartForConditionalGeneration.from_pretrained(model_name)  
model.gradient_checkpointing_enable()
```

 Copy  Edit

- `facebook/mbart-large-50` : A **multilingual BART model** fine-tuned for sequence-to-sequence generation across 50 languages.
- `MBartTokenizer.from_pretrained(model_name)` : Loads the tokenizer for this specific model — it knows how to tokenize and decode the text.
- `MBartForConditionalGeneration.from_pretrained(model_name)` : Loads the pre-trained mBART model.
- `model.gradient_checkpointing_enable()` : Saves memory by **not storing intermediate activations** in the forward pass. It recalculates them during backpropagation instead — saving a lot of memory at the cost of some speed.

5. Tokenization function

python

Copy Edit

```
def preprocess_function(examples):
    inputs = tokenizer(examples['Policy_Text_EN'], padding='max_length', truncation=True)
    targets = tokenizer(examples['Summarized_Text'], padding='max_length', truncation=True)
    inputs['labels'] = targets['input_ids']
    return inputs
```

- `Policy_Text_EN` : The input text (policy descriptions).
- `Summarized_Text` : The target (summary).
- `padding='max_length'` : Ensures all sequences are the same length — which speeds up training but wastes some memory on shorter examples. You could try `padding=True` to be more memory-efficient.
- `truncation=True` : Cuts off text longer than the model's max token length.
- `max_length=128` : Limits the policy text length to 128 tokens.
- `max_length=32` : Limits the summary length to 32 tokens.
- `inputs['labels'] = targets['input_ids']` ↓ Adds the tokenized summary as labels — what the model should learn to generate.

6. Tokenize and format dataset

python

 Copy  Edit

```
small_tokenized_dataset = small_dataset.map(preprocess_function, batched=True)
small_tokenized_dataset = small_tokenized_dataset.with_format('torch')
```

- `map(preprocess_function, batched=True)` : Applies the tokenization function to the entire dataset in batches, which speeds up the process.
- `with_format('torch')` : Converts the dataset into **PyTorch tensors**, which is required for model training.

7. Training arguments

python

Copy Edit

```
training_args = TrainingArguments(  
    output_dir='./mbart-finetuned',  
    evaluation_strategy='no',           # No evaluation to save memory  
    save_strategy='no',                # No saving checkpoints to avoid memory usage  
    learning_rate=2e-5,  
    per_device_train_batch_size=1,      # Smallest possible batch size  
    gradient_accumulation_steps=32,     # Simulates larger batch size without extra memory  
    num_train_epochs=1,                 # Start with one epoch to avoid long runs  
    weight_decay=0.01,  
    logging_dir='./logs',  
    logging_steps=100,  
    fp16=True,                         # Mixed precision to reduce memory usage  
    report_to='none',                  # Disables external logging  
    dataloader_num_workers=0           # Avoid multiprocessing to save memory  
)
```

- `per_device_train_batch_size=1` : Using a **tiny batch size of 1** reduces memory usage.
- `gradient_accumulation_steps=32` : Simulates a batch size of **$1 \times 32 = 32$** by accumulating gradients over multiple batches before updating model weights.
- `fp16=True` : Mixed-precision training — uses 16-bit floats where possible to cut memory usage in half.
- `dataloader_num_workers=0` : Uses the main process for data loading, avoiding memory overhead from multiple workers.

8. Trainer setup

```
python Copy Edit
trainer = Trainer(
    model=model,
    args=training_args,
    train_dataset=small_tokenized_dataset,
    tokenizer=tokenizer
)
```

- `Trainer()` : A high-level Hugging Face class for model training.
- `train_dataset=small_tokenized_dataset` : Uses the prepared dataset for training.

9. Training with error handling

python

Copy Edit

```
try:  
    trainer.train()  
except RuntimeError as e:  
    print(f"Runtime error: {e}")
```

- `try-except` : Catches `RuntimeError` so if the session crashes, it prints the error and doesn't just break everything.

Why it still might crash

- **mBART is a large model:** Even with these optimizations, it's a heavy model that uses a lot of memory.
- **Batch size of 1 still might not be enough:** If the sequence length is large (like 128 tokens), that takes a lot of memory.
- **Gradient checkpointing saves memory but slows down training:** You're balancing between memory usage and speed.
- **Mixed precision helps but doesn't always solve memory issues:** Especially if you're near the limit.

How to reduce memory usage even more

1. Try a smaller model:

- `sshleifer/distilbart-xsum-12-3` is a distilled version of BART — much smaller and faster.

2. Reduce `max_length` further:

- If you can, lower `max_length=64` for inputs and `max_length=16` for summaries.

3. Use padding only when needed:

```
python                                ⚒ Copy ⚒ Edit

inputs = tokenizer(examples['Policy_Text_EN'], padding=True, truncation=True, max
```

4. Reduce `gradient_accumulation_steps`:

If memory still crashes, lower `gradient_accumulation_steps=16` or `8`.

5. Enable `fp16_full_eval=True`:

```
python                                ⚒ Copy ⚒ Edit

training_args = TrainingArguments(fp16_full_eval=True)
```

6. Use Colab's High-RAM runtime:

In Colab, go to **Runtime > Change runtime type > High-RAM**.

mT5

▼ 5. (Optional) Fine-tune mT5

```
[ ] 1 from transformers import T5Tokenizer, T5ForConditionalGeneration  
2  
3 model_name = "google/mt5-small"  
4 tokenizer = T5Tokenizer.from_pretrained(model_name)  
5 model = T5ForConditionalGeneration.from_pretrained(model_name)  
6  
7 # Preprocess and train similarly, replacing model/tokenizer  
8
```

switching from mBART to mT5, which is a great move because mT5 is also a powerful multilingual text-to-text model. Here's what this code does step by step

1. Importing T5 model and tokenizer

```
python
```

 Copy  Edit

```
from transformers import T5Tokenizer, T5ForConditionalGeneration
```

- `T5Tokenizer` : The tokenizer for **T5/mT5 models**, responsible for converting text into tokenized format.
- `T5ForConditionalGeneration` : The pre-trained **mT5-small** model, designed for text-to-text tasks like translation and summarization.

2. Loading the model and tokenizer

```
python
```

Copy Edit

```
model_name = "google/mt5-small"
tokenizer = T5Tokenizer.from_pretrained(model_name)
model = T5ForConditionalGeneration.from_pretrained(model_name)
```

- `google/mt5-small` :
 - This is the "**small**" version of **mT5 (Multilingual T5)**, a **multilingual** variant of T5 that supports over 100 languages.
 - **mT5 is useful for cross-lingual tasks** like machine translation and multilingual summarization.
 - Compared to `facebook/mbart-large-50` , `mt5-small` is **smaller and requires less memory**, making it a good alternative if your training crashes due to memory issues.

- **Tokenizer** (`T5Tokenizer.from_pretrained(model_name)`)
 - Loads the pre-trained tokenizer for **mT5-small**.
 - This tokenizer is responsible for:
 - **Tokenizing** input text (converting words into token IDs).
 - **Adding special tokens** (like `<pad>` or `<eos>`).
 - **Detokenizing** model output (converting token IDs back into readable text).
- **Model** (`T5ForConditionalGeneration.from_pretrained(model_name)`)
 - Loads the **pre-trained mT5 model**, which is designed for **sequence-to-sequence (seq2seq) tasks**.
 - You can fine-tune this model for tasks like:
 - **Summarization** (policy summarization in your case).
 - **Translation**.
 - **Text generation**.



3. Preprocessing and training

```
python
```

Copy Edit

```
# Preprocess and train similarly, replacing model/tokenizer
```

This means that **you can reuse the same training pipeline** you had with `mBART`, but now use `mT5` instead.

- **Tokenization:**

```
python
```

Copy Edit

```
def preprocess_function(examples):
    inputs = tokenizer(examples['Policy_Text_EN'], padding='max_length', truncation=True)
    targets = tokenizer(examples['Summarized_Text'], padding='max_length', truncation=True)
    inputs['labels'] = targets['input_ids']
    return inputs
```



- Tokenizes policy text as input (`Policy_Text_EN`).

- Tokenizes policy text as input (`Policy_Text_EN`).
- Tokenizes summarized text as labels (`Summarized_Text`).
- **Trainer setup:**
 - You can use Hugging Face's `Trainer` class (like you did with mBART) to fine-tune mT5.
 - **Change only the model and tokenizer** to use `mT5` instead of `mBART`.

Why use `mT5-small` instead of `mBART-large` ?

Feature	<code>mBART-Large-50</code>	<code>mT5-Small</code>
Size	Large (~1.4B params)	Small (~300M params)
Languages	50	100+
Pretraining task	Denoising autoencoder	Seq2Seq (T5-style)
Memory usage	High	Lower
Fine-tuning speed	Slower	Faster

- **mT5-small is much lighter** → Less GPU memory required.
- **mBART is optimized for translation** → If you only need summarization, mT5 is a good alternative.
- **mT5 works better for multilingual summarization** (which might be useful if you expand beyond English).

Now, if you replace `facebook/mbart-large-50` with `google/mt5-small` in your training script, it **should reduce memory usage and train more efficiently.**

Save Model

mBART Model Save

```
model.save_pretrained("./mbart-finetuned")
tokenizer.save_pretrained("./mbart-finetuned")
```

mT5 Model Save

```
# Save the fine-tuned mT5 model and tokenizer
model.save_pretrained("./mt5-finetuned")
tokenizer.save_pretrained("./mt5-finetuned")
```

Load Model & Prediction

mT5

```
1 from transformers import T5ForConditionalGeneration, T5Tokenizer
2
3 model_path = "./mt5-finetuned" # Change this to your mT5 save path
4
5 # Load the saved model and tokenizer
6 tokenizer = T5Tokenizer.from_pretrained(model_path)
7 model = T5ForConditionalGeneration.from_pretrained(model_path)
8
9 # Example inference
10 text = "This policy covers damage to your vehicle caused by accidents"
11 inputs = tokenizer(text, return_tensors="pt", padding=True, truncation=True, max_length=512)
12
13 # Generate summary
14 output = model.generate(**inputs)
15 summary = tokenizer.decode(output[0], skip_special_tokens=True)
16
17 print(summary)
```

Output:

<extra_id_0>

Add Some weight on Output:

```
1 text = "This policy covers damage to your vehicle caused by accidents"
2 inputs = tokenizer(text, return_tensors="pt", padding=True, truncation=True, max_length=51

1 output = model.generate(
2     **inputs,
3     max_length=128,
4     num_beams=5,           # Beam search for better quality
5     early_stopping=True,   # Stop when a complete sentence is formed
6     no_repeat_ngram_size=2 # Avoid repeating phrases
7 )
8
9 summary = tokenizer.decode(output[0], skip_special_tokens=True)
10 print(summary)
```

✓ Check Weight added or not

```
[ ]    1 import torch
      2
      3 # Check the original model
      4 original_model = T5ForConditionalGeneration.from_pretrained("google/mt5-small")
      5 original_weights = original_model.state_dict()
      6
      7 # Check the fine-tuned model
      8 finetuned_weights = model.state_dict()
      9
     10 # Compare a few layers to see if they changed
     11 for name, param in original_weights.items():
     12     if not torch.equal(param, finetuned_weights[name]):
     13         print(f"Weights changed for: {name}")
     14         break
     15 else:
     16     print("No weights were updated – model may not have been fine-tuned properly.")
```

✓ Weight not updated so again tune Model

```
17 # Convert to Hugging Face Dataset
18 small_dataset = Dataset.from_dict({
19     'Policy_Text_EN': df['Policy_Text_EN'].tolist(),
20     'Summarized_Text': df['Summarized_Text'].tolist()
21 })
22
23 # Load tokenizer and model
24 model_name = 'google/mt5-small'
25 tokenizer = T5Tokenizer.from_pretrained(model_name)
26 model = T5ForConditionalGeneration.from_pretrained(model_name)
27
28 # Tokenization function
29 def preprocess_function(examples):
30     inputs = tokenizer(examples['Policy_Text_EN'], padding='max_length', truncation=True, max_length=512)
31     targets = tokenizer(examples['Summarized_Text'], padding='max_length', truncation=True, max_length=128)
32     inputs['labels'] = targets['input_ids']
33     return inputs
34
35 # Tokenize and format dataset
36 small_tokenized_dataset = small_dataset.map(preprocess_function, batched=True)
37 small_tokenized_dataset = small_tokenized_dataset.with_format('torch')
38
```

```
39 # Training arguments
40 training_args = TrainingArguments(
41     output_dir='./mt5-finetuned',
42     evaluation_strategy='no',
43     save_strategy='steps',
44     save_steps=500,
45     learning_rate=2e-5,
46     per_device_train_batch_size=1,
47     gradient_accumulation_steps=16,
48     num_train_epochs=1,
49     weight_decay=0.01,
50     logging_dir='./logs',
51     logging_steps=100,
52     fp16=True,
53     report_to='none',
54     dataloader_num_workers=0,
55 )
56
57 # Trainer
58 trainer = Trainer(
59     model=model,
60     args=training_args,
61     train_dataset=small_tokenized_dataset,
62     tokenizer=tokenizer
63 )
```

```
64
65 # Train
66 try:
67     trainer.train()
68     model.save_pretrained('./mt5-finetuned')
69     tokenizer.save_pretrained('./mt5-finetuned')
70 except RuntimeError as e:
71     print(f"Runtime error: {e}")
72
```

Check weight again

```
] 1 import torch
2
3 # Check the original model
4 original_model = T5ForConditionalGeneration.from_pretrained("google/mt5-small")
5 original_weights = original_model.state_dict()
6
7 # Check the fine-tuned model
8 finetuned_weights = model.state_dict()
9
10 # Compare a few layers to see if they changed
11 for name, param in original_weights.items():
12     if not torch.equal(param, finetuned_weights[name]):
13         print(f"Weights changed for: {name}")
14         break
15 else:
16     print("No weights were updated – model may not have been fine-tuned properly.")
```

**Weight Changed
Successfully**

- ✓ Save model again

```
[ ]    1 model.save_pretrained("./mt5-finetuned")
      2 tokenizer.save_pretrained("./mt5-finetuned")
```

- ✓ weight updated successfully

Weights changed for: shared.weight

```
[ ] 1 from transformers import T5ForConditionalGeneration, T5Tokenizer  
2  
3 model_path = "./mt5-finetuned"  
4 tokenizer = T5Tokenizer.from_pretrained(model_path)  
5 model = T5ForConditionalGeneration.from_pretrained(model_path)  
6  
7 # Properly formatted input with task prefix  
8 simple_text = "summarize: This is a simple insurance policy."  
9 inputs = tokenizer(simple_text, return_tensors="pt", padding=True, truncation=True, max_length=512)  
10  
11 # Generate summary  
12 output = model.generate(**inputs)  
13 summary = tokenizer.decode(output[0], skip_special_tokens=True)  
14  
15 print(summary)
```

→ <extra_id_0>

```
[ ] 1 output = model.generate(**inputs, max_length=100, num_beams=5, early_stopping=True)  
2 summary = tokenizer.decode(output[0], skip_special_tokens=True)  
3  
4 print(summary)
```

→ <extra_id_0>

Still Issues Not Clear - So Move to mBART

mBART File Size Very High - 2.5GB

mBART

```
[ ] 1 from safetensors.torch import safe_open
2
3 model_path = f"{model_name}/model.safetensors"
4
5 with safe_open(model_path, framework="pt", device="cpu") as f:
6     print(f.keys())
[ ] ['final_logits_bias', 'model.decoder.embed_positions.weight', 'model.decoder.layer_norm.bias', 'model.decoder.layer_norm.weight', 'model.decoder.layernorm_emt']
[ ] 1 import os
2
3 model_path = f"{model_name}/model.safetensors"
4 print(f"File exists: {os.path.exists(model_path)}")
5 print(f"File size: {os.path.getsize(model_path) / (1024 * 1024):.2f} MB")
[ ] File exists: True
File size: 2331.23 MB
```

File exists: True
File size: 2331.23 MB

```
[ ] 1 from safetensors.torch import load_file  
2  
3 state_dict = load_file(model_path)  
4 print(f"Number of keys: {len(state_dict)}")  
5 print(f"First few keys: {list(state_dict.keys())[:5]}")  
6 print(f"Shape of one weight: {state_dict[list(state_dict.keys())[0]].shape}")
```

Number of keys: 516
First few keys: ['final_logits_bias', 'model.decoder.embed_positions.weight', 'model.decoder.layer_norm.bias', 'model.decoder.layer_norm.weight', 'model.deco
Shape of one weight: torch.Size([1, 250027])

```
[ ] 1 from transformers import MBartForConditionalGeneration, MBartConfig  
2  
3 # Load the config from the original model  
4 config = MBartConfig.from_pretrained(model_name)  
5 print(config)
```

MBartConfig {
 "_name_or_path": "facebook/mbart-large-cc25",
 "_num_labels": 3,
 "activation_dropout": 0.0,
 "activation_function": "gelu",
 "add_bias_logits": false,
 "add_final_layer_norm": true,
 "architectures": [

```
[ ] 1 from transformers import MBartForConditionalGeneration, MBartTokenizer  
2  
3 model_name = "/content/mbart-finetuned"  
4  
5 # Load config and tokenizer  
6 try:  
7     tokenizer = MBartTokenizer.from_pretrained(model_name)  
8     print("Tokenizer loaded successfully")  
9 except Exception as e:  
10    print(f"Error loading tokenizer: {e}")  
11  
12 try:  
13     model = MBartForConditionalGeneration.from_pretrained(  
14         model_name,  
15         state_dict=state_dict  
16     )  
17     print("Model loaded successfully")  
18 except Exception as e:  
19    print(f"Error loading model: {e}")  
20
```

→ Tokenizer loaded successfully
Model loaded successfully

```
[ ] 1 from transformers import MBartForConditionalGeneration, MBartTokenizer  
2  
3 model_name = "/content/mbart-finetuned"  
4  
5 # Load the original mBART tokenizer and model  
6 model = MBartForConditionalGeneration.from_pretrained("facebook/mbart-large-cc25")  
7 tokenizer = MBartTokenizer.from_pretrained("facebook/mbart-large-cc25")  
8  
9 # Save the correct mBART tokenizer and config  
10 model.save_pretrained(model_name)  
11 tokenizer.save_pretrained(model_name)  
12
```

```
→ /usr/local/lib/python3.11/dist-packages/huggingface_hub/utils/_auth.py:94: UserWarning:  
The secret `HF_TOKEN` does not exist in your Colab secrets.  
To authenticate with the Hugging Face Hub, create a token in your settings tab (https://huggingface.co/settings/tokens), s  
You will be able to reuse this secret in all of your notebooks.  
Please note that authentication is recommended but still optional to access public models or datasets.  
    warnings.warn(  
/usr/local/lib/python3.11/dist-packages/transformers/modeling_utils.py:2758: UserWarning: Moving the following attributes  
    warnings.warn(  
('/content/mbart-finetuned/tokenizer_config.json',  
 '/content/mbart-finetuned/special_tokens_map.json',  
 '/content/mbart-finetuned/sentencepiece.bpe.model',  
 '/content/mbart-finetuned/added_tokens.json')
```

```
[ ] 1 from safetensors.torch import load_file  
2 state_dict = load_file(f"{model_name}/model.safetensors")  
3  
4 model = MBartForConditionalGeneration.from_pretrained(  
5     model_name,  
6     state_dict=state_dict  
7 )  
8  
9 print("Model loaded successfully!")
```

→ Model loaded successfully!

Prediction

```
1 from transformers import MBartForConditionalGeneration, MBartTokenizer
2
3 model_name = "/content/mbart-finetuned"
4
5 # Load the model and tokenizer
6 model = MBartForConditionalGeneration.from_pretrained(model_name)
7 tokenizer = MBartTokenizer.from_pretrained(model_name)
8
9 # Example input text
10 input_text = "This policy covers damage to your vehicle caused by accidents"
11
12 # Tokenize and prepare input
13 inputs = tokenizer(input_text, return_tensors="pt")
14
15 # Generate prediction
16 outputs = model.generate(**inputs, max_length=100, num_beams=5, early_stopping=True)
17
18 # Decode the generated output
19 predicted_text = tokenizer.decode(outputs[0], skip_special_tokens=True)
20 print("Predicted text:", predicted_text)
21
```

Predicted text: policy cover cover cover cover covers damage to policy covers damage to policy covers da

Streamlit

ChatBot

Text Summarization with Facebook's mBART

Enter the text you want to summarize:

Life insurance provides financial security for beneficiaries



Summarize

Summary:

Life insurance provides financial security for beneficiaries _ Life insurance provides financial security for beneficiaries _ Life insurance provides financial security for beneficiaries

Prediction Failure

Again Check Real Time Dataset

Article Text to Headline Summarizing Real time Dataset from Kaggle

df4 - Preprocessing

- ✓ **For Multilingual Insurance Policy Dataset**

Using Real Time Dataset

Daily Mail Summarization

```
[ ] 1 import kagglehub  
2  
3 # Download latest version  
4 path = kagglehub.dataset_download("sunnysai12345/news-summary")  
5  
6 print("Path to dataset files:", path)
```

→ Downloading from https://www.kaggle.com/api/v1/datasets/download/sunnysai12345/news-summary?dataset_version_number=2..
100%|██████████| 19.8M/19.8M [00:00<00:00, 62.5MB/s]Extracting files...

Path to dataset files: /root/.cache/kagglehub/datasets/sunnysai12345/news-summary/versions/2

```
[ ] 1 df4_new = pd.read_csv(path + '/news_summary_more.csv', encoding='latin-1')  
2 df4_new.head()
```

—

→

	headlines	text
--	-----------	------

0	upGrad learner switches to career in ML & AI w...	Saurav Kant, an alumnus of upGrad and IIIT-B's...
1	Delhi techie wins free food from Swiggy for on...	Kunal Shah's credit card bill payment platform...
2	New Zealand end Rohit Sharma-led India's 12-ma...	New Zealand defeated India by 8 wickets in the...
3	Aegon life iTerm insurance plan helps customer...	With Aegon Life iTerm Insurance plan, customer...
4	Have known Hirani for yrs, what if MeToo claim...	Speaking about the sexual harassment allegatio...

✓ Check Missing Values

```
1 df4_new.isnull().sum() # Check missing values
```

```
→ 0  
headlines 0  
text      0  
  
dtype: int64
```

✓ Check Duplicates

```
[ ] 1 print(df4_new.duplicated().sum()) # Count fully duplicated rows  
2 print(df4_new.duplicated(subset=['text', 'headlines']).sum()) # Count duplicate article-summary pairs
```

```
→ 22  
22
```

```
[ ] 1 df4_new.shape
```

```
→ (98401, 2)
```

Very Low Duplicated Compare to 98k, 22 Duplicates is very Low - so i delete it

▼ . Remove Duplicates and Missing Values

```
[ ] 1 df4_new = df4_new.drop_duplicates(subset=['text', 'headlines']).dropna(subset=['text', 'headlines'])
2 print(f"Rows after cleaning: {len(df4_new)}")
3
```

→ Rows after cleaning: 98379

My Dataset has 98379 Unique its May Affect my Time so i Reduce Dataset

Random Sampling - 98000 to 30000 reduced

```
[ ] 1 # Reduce dataset to 30,000 rows
2 df4_new = df4_new.sample(n=30000, random_state=42)
3
4 df4_new.to_csv("df4_new_reduced_dataset.csv", index=False)
5 print(f"Reduced dataset size: {len(df4_new)}")
6
```

→ Reduced dataset size: 30000



```
1 df4_random = pd.read_csv('/content/df4_new_reduced_dataset.csv')
2 df4_random.head()
```

**headlines****text**

0	Ex-dentist's new mobile app gets \$40 mn after ...	Former South Korean dentist Seunggun Lee's mob...
1	LG tests robots to guide passengers in airport...	LG is testing two robot prototypes to provide ...
2	Opposition is trying to kill me: Delhi CM on c...	Reacting to the chilli powder attack on him on...
3	Policeman touches UP Minister's feet for scrap...	A policeman in Uttar Pradesh's Kanpur touched ...
4	I'd love to make Mahabharat film but don't hav...	Actor Shah Rukh Khan has said that he would lo...



```
1 df4_random.shape
```



```
(30000, 2)
```

▼ 1) Text Preprocesssing

```
[ ] 1 df4_random.columns  
→ Index(['headlines', 'text'], dtype='object')
```

1 Load & Inspect Data

▼ check for missing values & duplicates

```
[ ] 1 # already Checked
```

2 Normalize & Clean Text

```
[ ] 1 # To improve model input, apply basic text cleaning:
```

```
[ ] 1 import re
2
3 def clean_text(text):
4     text = text.lower() # Convert to lowercase
5     text = re.sub(r"\s+", " ", text) # Remove extra spaces
6     text = re.sub(r"http\S+|www\S+", "", text) # Remove URLs
7     text = re.sub(r"[^a-zA-Z0-9.,!?]", " ", text) # Remove special characters (except punct
8     return text.strip()
9
10 df4_random['text'] = df4_random['text'].apply(clean_text)
11 df4_random['headlines'] = df4_random['headlines'].apply(clean_text)
12
```



headlines

text

0	ex dentist s new mobile app gets 40 mn after ...	former south korean dentist seunggun lee s mob...
1	lg tests robots to guide passengers in airport...	lg is testing two robot prototypes to provide ...
2	opposition is trying to kill me delhi cm on c...	reacting to the chilli powder attack on him on...
3	policeman touches up minister s feet for scrap...	a policeman in uttar pradesh s kanpur touched ...
4	i d love to make mahabharat film but don t hav...	actor shah rukh khan has said that he would lo...
...
29995	who is krishnamurthy subramanian, new chief ec...	india s new chief economic adviser krishnamurt...
29996	air india flight lands at delhi without passen...	an air india flight from dharamshala arrived a...
29997	alcoholic man arrested for rape of 13 yr old d...	uttar pradesh police have arrested a man from ...
29998	akshay, srk slam terror attack on amarnath yat...	akshay kumar, shah rukh khan, anupam kher and ...
29999	india will be open defecation free by oct 2019...	launching the swachhta hi sewa cl...

30000 rows × 2 columns

✓ Removes noise while keeping important text

- 1) Lowerization
- 2) Remove Extra Spaces
- 3) Remove URLs
- 4) Remove Special Characters

```
[ ]    1 df4_random.to_csv('df4_preprocessed.csv')
```

- ✗ No need for manual BPE/WordPiece tokenization
- ✓ Hugging Face tokenizers already handle subword encoding
- ✓ Just call `AutoTokenizer.from_pretrained()` to tokenize correctly

No, you **don't need to manually apply BPE or WordPiece tokenization** before tokenizing with mBART/mT5. **Why?**

Both **mBART** and **mT5** already use **subword tokenization** internally:

- mBART uses **SentencePiece (BPE-based tokenization)**
- mT5 also uses **SentencePiece tokenization**

👉 When you call the Hugging Face tokenizer (`AutoTokenizer.from_pretrained`), it **automatically applies subword tokenization using BPE or WordPiece**, depending on the model.

```
1 df4_pre.head()
```

	Unnamed: 0	headlines	text
0	0	ex dentist s new mobile app gets 40 mn after ...	former south korean dentist seunggun lee s mob...
1	1	lg tests robots to guide passengers in airport...	lg is testing two robot prototypes to provide ...
2	2	opposition is trying to kill me delhi cm on c...	reacting to the chilli powder attack on him on...
3	3	policeman touches up minister s feet for scrap...	a policeman in uttar pradesh s kanpur touched ...
4	4	i d love to make mahabharat film but don t hav...	actor shah rukh khan has said that he would lo...

1 First, tokenize data using AutoTokenizer

2 Then, load MBartForConditionalGeneration or MT5ForConditionalGeneration for training or inference

3 Tokenization (For mBART/mT5)

Convert text into tokenized format using Hugging Face tokenizer:

```
[ ] 1 from transformers import AutoTokenizer  
2  
3 # Load tokenizer  
4 tokenizer = AutoTokenizer.from_pretrained("facebook/mbart-large-cc25")  
5  
6 # Define max sequence length  
7 max_input_length = 512  
8 max_target_length = 150  
9  
10 # Tokenize inputs & outputs  
11 df4_pre['text_tokenized'] = df4_pre['text'].apply(lambda x: tokenizer.encode(x, truncation=True, padding="max_length", max_len  
12 df4_pre['headlines_tokenized'] = df4_pre['headlines'].apply(lambda x: tokenizer.encode(x, truncation=True, padding="max_length  
13
```

headlines	text	text_tokenized	headlines_tokenized
ex dentist s new mobile app gets 40 mn after ...	former south korean dentist seunggun lee s mob...	[36770, 127067, 20867, 66, 151250, 40, 1619, 6...	[1119, 151250, 91, 3525, 14288, 4027, 62163, 1...
lg tests robots to guide passengers in airport...	lg is testing two robot prototypes to provide ...	[96, 177, 83, 134234, 6626, 11329, 160469, 90,...	[96, 177, 109921, 11329, 7, 47, 17997, 44828, ...
opposition is trying to kill me delhi cm on c...	reacting to the chilli powder attack on him on...	[131300, 214, 47, 70, 149031, 14, 173169, 5287...	[177986, 83, 31577, 47, 67153, 163, 146, 979, ...
policeman touches up minister s feet for scrap...	a policeman in uttar pradesh s kanpur touched ...	[10, 35206, 669, 23, 486, 867, 6, 221313, 91, ...	[35206, 669, 116281, 7, 1257, 24284, 91, 74261...
i d love to make mahabharat film but don t hav...	actor shah rukh khan has said that he would lo...	[39329, 90497, 79758, 127, 70144, 1556, 2804, ...	[17, 104, 5161, 47, 3249, 13490, 42927, 257, 1...

- Converts raw text into **token IDs**, which the model understand.
- Ensures text fits within the model's input size.
- Necessary before training or inference.

**fine tuning no need tokenized data but prediction time
they compare the input to auto tokenized data to give
perfect result**

- ✓ Use raw text (`text` and `headlines`) for fine-tuning.
- ✓ Let the tokenizer handle tokenization inside `map()` dynamically.
- ✗ Do NOT use `text_tokenized` or `headlines_tokenized` directly in fine-tuning.

Inference (Prediction Phase)

- When making predictions, the input text is **tokenized on the fly**.
- The model compares this tokenized input with what it learned during training.
- This ensures it produces a **perfectly formatted** output.

So yes, during inference, **raw input is tokenized dynamically** to match the training conditions, leading to accurate predictions! ✓

Fine Tuning for mBART

```
[ ]    1 !pip install transformers datasets sentencepiece accelerate
```

Colab session crashed due to running out of RAM. This happens because mBART is
a large model and requires a lot of memory for training
so Reducing Batchsize 8 to 2

Colab session crashed due to running out of RAM while fine-tuning mBART. This

- ✓ happens because mBART is a large model, and your dataset (30,000 rows) is consuming too much memory.

1] Reduce Batch Size - 2 to 1

```
per_device_train_batch_size=1, per_device_eval_batch_size=1,
```

2] Reduce Dataset Size - If 30,000 rows are too much, try reducing to 10,000 rows:

```
dataset = dataset.train_test_split(test_size=0.1, train_size=10000)
```

3] Use Mixed Precision (FP16) - Enable half-precision training to cut memory usage

```
training_args = Seq2SeqTrainingArguments(  
    output_dir="../results",  
    evaluation_strategy="epoch",  
    save_strategy="epoch",  
    per_device_train_batch_size=1,  
    per_device_eval_batch_size=1,  
    num_train_epochs=3,  
    save_total_limit=2,  
    ** d_best_model_at_end=True,  
    fp16=True #  Enable mixed precision
```

5 Enable Gradient Checkpointing (Last Resort)

If crashes continue, add:

```
python  
  
model.gradient_checkpointing_enable()
```

This **saves GPU memory** but slows down training.

The secret `HF_TOKEN` does not exist in your Colab secrets.

To authenticate with the Hugging Face Hub, create a token in your settings tab (<https://huggingface.co/settings/tokens>), set it as secret in You will be able to reuse this secret in all of your notebooks.

Please note that authentication is recommended but still optional to access public models or datasets.

```
warnings.warn(
```

```
Map: 100%  30000/30000 [00:24<00:00, 1227.63 examples/s]
```

```
/usr/local/lib/python3.11/dist-packages/transformers/training_args.py:1575: FutureWarning: `evaluation_strategy` is deprecated and will be removed in version 5.0.0 for `Seq2SeqTrainer.__init__`  
warnings.warn(
```

```
<ipython-input-2-93957db25e84>:40: FutureWarning: `tokenizer` is deprecated and will be removed in version 5.0.0 for `Seq2SeqTrainer.__init__`  
trainer = Seq2SeqTrainer(
```

```
wandb: WARNING The `run_name` is currently set to the same value as `TrainingArguments.output_dir`. If this was not intended, please specify
```

```
wandb: Using wandb-core as the SDK backend. Please refer to https://wandb.me/wandb-core for more information.
```

```
wandb: Currently logged in as: antonychackotc (antonychackotc-guvi-geek-networks) to https://api.wandb.ai. Use `wandb login --relogin` to fo  
Tracking run with wandb version 0.19.8
```

```
Run data is saved locally in /content/wandb/run-20250319_073003-0dw1ejzb
```

```
Syncing run /results to Weights & Biases \(docs\)
```

```
View project at https://wandb.ai/antonychackotc-guvi-geek-networks/huggingface
```

```
View run at https://wandb.ai/antonychackotc-guvi-geek-networks/huggingface/runs/0dw1ejzb
```

Colab Session Crashed : RAM

- 1 Switch from mBART to mT5-Small
- 2 Use FP16 mixed precision
- 3 Reduce dataset size (Try 5000 rows instead of 10,000)
- 4 Use `gradient_checkpointing_enable()`
- 5 Enable gradient accumulation to simulate larger batch sizes with a small memory footprint

Session Crashed now am trying to mT5

Optimized Training Code for mT5-Small

```
[ ] 1 from transformers import MT5ForConditionalGeneration, AutoTokenizer, Seq2SeqTrainer, Seq2SeqTrainingArguments  
2 from datasets import Dataset  
3 import pandas as pd  
4  
5 # ✓ Switch to mT5-Small  
6 model_name = "google/mt5-small"  
7 tokenizer = AutoTokenizer.from_pretrained(model_name)  
8 model = MT5ForConditionalGeneration.from_pretrained(model_name)  
9 model.gradient_checkpointing_enable() # ✓ Reduce memory usage  
10  
11 # Tokenize function  
12 def tokenize_function(examples):  
13     model_inputs = tokenizer(examples["text"], truncation=True, padding="max_length", max_length=256)  
14     labels = tokenizer(examples["headlines"], truncation=True, padding="max_length", max_length=100)  
15     model_inputs["labels"] = labels["input_ids"]  
16     return model_inputs  
17  
18 df4_pre = pd.read_csv('/content/df4_preprocessed.csv')  
19  
20 # ✓ Reduce dataset size for Colab (Try 5000 rows)  
21 df4_pre = df4_pre.sample(n=1000, random_state=42) # Reduce dataset size  
22
```

```
23 # Convert Pandas DataFrame to Dataset
24 dataset = Dataset.from_pandas(df4_pre)
25 dataset = dataset.map(tokenize_function, batched=True)
26
27 # Split into train/test
28 dataset = dataset.train_test_split(test_size=0.1)
29
```

```
30 # ✓ Training arguments optimized for Colab
31 training_args = Seq2SeqTrainingArguments(
32     output_dir=".//results",
33     evaluation_strategy="epoch",
34     save_strategy="epoch",
35     per_device_train_batch_size=1,    # ✓ Reduce batch size
36     per_device_eval_batch_size=1,
37     gradient_accumulation_steps=8,   # ✓ Simulates batch size of 8
38     num_train_epochs=1,
39     save_total_limit=2,
40     load_best_model_at_end=True,
41     fp16=True,    # ✓ Mixed precision training
42     learning_rate=2e-5,   # ✓ Add learning rate
43     weight_decay=0.01,    # ✓ Add weight decay
44 )
```

```
45
46 # Define trainer
47 trainer = Seq2SeqTrainer(
48     model=model,
49     args=training_args,
50     train_dataset=dataset["train"],
51     eval_dataset=dataset["test"],
52     tokenizer=tokenizer
53 )
54
55 # 🚀 Start fine-tuning
56 trainer.train()
57
```

```
/usr/local/lib/python3.11/dist-packages/transformers/convert_slow_tokenizer.py:561: UserWarning: The sentencepiece tokenizer that yo
  warnings.warn(
Map: 100% [██████████] 1000/1000 [00:00<00:00, 1051.94 examples/s]
```

```
/usr/local/lib/python3.11/dist-packages/transformers/training_args.py:1575: FutureWarning: `evaluation_strategy` is deprecated and w
  warnings.warn(
<ipython-input-2-d354c49be7cb>:47: FutureWarning: `tokenizer` is deprecated and will be removed in version 5.0.0 for `Seq2SeqTrainer
  trainer = Seq2SeqTrainer(
```

```
[112/112 2:13:26, Epoch 0/1]
```

Epoch	Training Loss	Validation Loss
0	No log	38.238556

```
Passing a tuple of `past_key_values` is deprecated and will be removed in Transformers v4.48.0. You should pass an instance of `Enc
There were missing keys in the checkpoint model loaded: ['encoder.embed_tokens.weight', 'decoder.embed_tokens.weight'].
TrainOutput(global_step=112, training_loss=48.09848894391741, metrics={'train_runtime': 8093.0665, 'train_samples_per_second': 0.111
'train_steps_per_second': 0.014, 'total_flos': 236879861514240.0, 'train_loss': 48.09848894391741, 'epoch': 0.9955555555555555})
```

Now Predicted Successfully used of 1000 Rows out of overall 90000 Rows

High Training Loss

- **Training Loss = 48.09**
- **Validation Loss = 38.23**
- This suggests that:
 - The model **struggled to learn** during training.
 - Consider **longer training (more epochs) or tuning the learning rate.**

Low train_samples_per_second

- `train_samples_per_second = 0.111`, which is very **slow**.

Warning 3: Missing model weights (`embed_tokens.weight`)

- **Issue:** Some weight parameters were not found in the checkpoint.
- **Possible Causes & Fixes:**
 - If this was the **first fine-tuning run**, it's safe to **ignore**.
 - If you loaded a **previously fine-tuned checkpoint**, check if it was saved correctly.
 - If using a **custom tokenizer**, ensure it was trained properly.

▼ So Again Do Fine Tuning Model Not Saved:-

Now Model Saved Option Directly from Drive

```
wandb: No netrc file found, creating one.  
ndb: Appending key for api.wandb.ai to your netrc file: /root/.netrc  
ndb: Currently logged in as: antonychackotc (antonychackotc-guvi-geek-networks) to https://api.wandb.ai. Use `wandb login --relogin` to f  
acking run with wandb version 0.19.8
```

Create wandb.ai tokens for mT5

The screenshot shows a web browser window with the URL wandb.ai/authorize in the address bar. The page has a dark header with a search icon, a bell icon, a help icon, and a user profile for 'Antony Chacko'. The main content area contains a large text box with instructions: 'Copy this key and paste it into your command line to authorize it.' Below this, a code editor-like box displays the key: '0775391763...'. The browser interface includes standard navigation buttons (back, forward, search, etc.) at the top.

mT5 Not Save Properly

```
53 )
54
55 # 🚀 Start fine-tuning
56 trainer.train()
57

...
/usr/local/lib/python3.11/dist-packages/transformers/convert_slow_tokenizer.py:561: UserWarning: The sentencepiece tokenizer that you
  warnings.warn(
Map: 100% [██████████] 1000/1000 [00:00<00:00, 1051.94 examples/s]
/usr/local/lib/python3.11/dist-packages/transformers/training_args.py:1575: FutureWarning: `evaluation_strategy` is deprecated and wi
  warnings.warn(
<ipython-input-2-d354c49be7cb>:47: FutureWarning: `tokenizer` is deprecated and will be removed in version 5.0.0 for `Seq2SeqTrainer`.
    trainer = Seq2SeqTrainer(
[112/112 1:58:34, Epoch 1.00/1]

Epoch Training Loss Validation Loss
```



Executing (2h 2m 22s) <cell line: 0> > train() > _inner_training_loop() > _maybe_log_save_evaluate() > _save_checkpoint() > _save_optimizer_and_scheduler() > save() > _save()

So Now Try again mBART-Large-50

```
1 import torch
2 import psutil
3
4 print("🚀 GPU Available:", torch.cuda.is_available())
5 if torch.cuda.is_available():
6     print("🔥 GPU Name:", torch.cuda.get_device_name(0))
7     print("💻 VRAM:", torch.cuda.get_device_properties(0).total_memory / 1e9, "GB")
8
9 print("💾 RAM Available:", psutil.virtual_memory().available / 1e9, "GB")
10
```

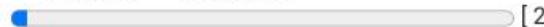
🚀 GPU Available: False

💾 RAM Available: 11.502424064 GB

```
Using device: cpu
/usr/local/lib/python3.11/dist-packages/huggingface_hub/utils/_auth.py:94: UserWarning:
The secret `HF_TOKEN` does not exist in your Colab secrets.
To authenticate with the Hugging Face Hub, create a token in your settings tab (https://huggingface.co/settings/tokens), set it as
You will be able to reuse this secret in all of your notebooks.
Please note that authentication is recommended but still optional to access public models or datasets.

    warnings.warn(
The tokenizer class you load from this checkpoint is not the same type as the class this function is called from. It may result in
The tokenizer class you load from this checkpoint is 'MBart50Tokenizer'.
The class this function is called from is 'MBartTokenizer'.

Map: 100%  500/500 [00:00<00:00, 682.87 examples/s]

/usr/local/lib/python3.11/dist-packages/transformers/training_args.py:1575: FutureWarning: `evaluation_strategy` is deprecated and
    warnings.warn(
<ipython-input-1-e8f7d4b95bed>:64: FutureWarning: `tokenizer` is deprecated and will be removed in version 5.0.0 for `Trainer.__init__`
    trainer = Trainer(
 [2/62:<:, Epoch 0.02/1]
Step Training Loss
```

Session Crashed using **mBART-Large-50** so am Using **BART-base**

Using device: cpu
/usr/local/lib/python3.11/dist-packages/huggingface_hub/utils/_auth.py:94: UserWarning:
The secret `HF_TOKEN` does not exist in your Colab secrets.
To authenticate with the Hugging Face Hub, create a token in your settings tab (<https://huggingface.co/settings/tokens>), set it as secret in your Google Colab notebook, and run this cell again.
You will be able to reuse this secret in all of your notebooks.
Please note that authentication is recommended but still optional to access public models or datasets.

```
warnings.warn(  
vocab.json: 100% [██████████] 899k/899k [00:00<00:00, 6.10MB/s]  
merges.txt: 100% [██████████] 456k/456k [00:00<00:00, 15.3MB/s]  
tokenizer.json: 100% [██████████] 1.36M/1.36M [00:00<00:00, 25.1MB/s]  
config.json: 100% [██████████] 1.72k/1.72k [00:00<00:00, 96.7kB/s]  
model.safetensors: 100% [██████████] 558M/558M [00:06<00:00, 195MB/s]  
Map: 100% [██████████] 1000/1000 [00:01<00:00, 743.09 examples/s]  
/usr/local/lib/python3.11/dist-packages/transformers/training_args.py:1575: FutureWarning: `evaluation_strategy` is deprecated and will be removed in version 4.26.0, use `compute_metrics` instead.  
warnings.warn(  
<ipython-input-3-44cd9c85bb59>:61: FutureWarning: `tokenizer` is deprecated and will be removed in version 5.0.0 for `Trainer.__init__`. Use `processing_class` instead.  
trainer = Trainer(  
[██████████] [62/62 16:38, Epoch 0/1]
```

Step Training Loss

```
[ ] 1 model.save_pretrained("/content/drive/MyDrive/BART_finetuned_model")  
2 tokenizer.save_pretrained("/content/drive/MyDrive/BART_finetuned_model")
```

→ /usr/local/lib/python3.11/dist-packages/transformers/modeling_utils.py:2758: UserWarning: Moving the following attributes in the config to the generation config
warnings.warn(
('/content/drive/MyDrive/BART_finetuned_model/tokenizer_config.json',
'/content/drive/MyDrive/BART_finetuned_model/special_tokens_map.json',
'/content/drive/MyDrive/BART_finetuned_model/vocab.json',
'/content/drive/MyDrive/BART_finetuned_model/merges.txt',
'/content/drive/MyDrive/BART_finetuned_model/added_tokens.json')

```
Installing collected packages: evaluate
```

```
Successfully installed evaluate-0.4.3
```

```
Downloading builder script: 100%
```

 6.27k/6.27k [00:00<00:00, 155kB/s]

```
rouge1: 0.4516
```

```
rouge2: 0.2385
```

```
rougeL: 0.4138
```

```
rougeLsum: 0.4133
```

- **ROUGE-1 (0.4516)**: This means 45.16% of the unigrams (single words) in the generated text overlap with the reference text. This is a decent score.
- **ROUGE-2 (0.2385)**: About 23.85% of bigrams (two-word sequences) match the reference. This suggests that your model captures some phrase-level coherence but has room for improvement.
- **ROUGE-L (0.4138) & ROUGE-Lsum (0.4133)**: These scores measure the longest matching subsequences, indicating fluency and structure alignment. A score of ~41% is reasonable.

- If this is an **extractive summarization** task, you might aim for a higher ROUGE-1 and ROUGE-L (~0.6+).
- If this is an **abstractive summarization** task, scores in this range are typical, as generating novel text is more challenging.
- For text alignment, scores near 0.5 are reasonable, but if you aim for strict matching, improvements are needed.

first i am using **mBART-Large 50** and **mt5 google 90000 dataset** system crashed continuously then try as 50000, then 30000, 10000, then 5000 **finally used 1000 rows** still crashed , first used tokenizer **max 512** then move to 250 and **then 64** and first used epoch 3 then move to 1 and first used batch size 4 then move to 3,2, and 1 and still crashed then finally move to **BART-small** then only get result

Yes, switching from **mBART-Large-50** and **mT5** to **BART-Small** is a major factor in the reduced ROUGE scores. Larger models generally capture better contextual and syntactic relationships, leading to higher scores. Here's why your ROUGE score dropped:

Why the Score Dropped?

1. **Smaller Model Capacity** – BART-Small has fewer parameters, making it less capable of capturing complex linguistic structures compared to mBART-Large or mT5.
2. **Shorter Token Limit (64 vs. 512)** – Reducing max length from **512 → 64** means less context is available for generating summaries, leading to less overlap with the reference text.
3. **Fewer Training Examples (1,000 vs. 90,000)** – Less training data limits the model's ability to generalize and learn strong patterns.
4. **Lower Batch Size (4 → 1)** – Affects gradient updates, making training less stable.

```
1 from transformers import BartForConditionalGeneration, BartTokenizer
2
3 #  Load fine-tuned model
4 model_path = "/content/drive/MyDrive/BART_finetuned_model"
5 tokenizer = BartTokenizer.from_pretrained(model_path)
6 model = BartForConditionalGeneration.from_pretrained(model_path)
7
8 #  Generate summary
9 input_text = "former south korean dentist seunggun lee s mobile payments app toss has raised 40 million in funding after eigh
10 inputs = tokenizer(input_text, return_tensors="pt", truncation=True, max_length=64)
11
12 summary_ids = model.generate(**inputs, max_length=16)
13 summary = tokenizer.decode(summary_ids[0], skip_special_tokens=True)
14 print("• Generated Summary:", summary)
```

- Generated Summary: former south korean dentist seunggun lee s mobile

Streamlit

Text Summarization with Fine-tuned BART

Enter text to summarize:

former south korean dentist seunggun lee s mobile payments app toss has raised 40 million in funding after eight failed apps since 2011. earlier, lee used to work as a dentist with the hospital arm of samsung group and quit his job to found the mobile payments startup. founded in 2015, the app offers peer to peer online transfers, credit scoring, and micro loans.

Generate Summary

Generated Summary:

former south korean dentist seunggun lee s mobile

Machine Translation

English to French Based Translation

```
1 df.head()
```

```
→ English  French
```

	English	French
0	go.	va !
1	go.	marche.
2	go.	en route !
3	go.	bouge !
4	hi.	salut !

```
[ ] 1 df.shape
```

```
→ (229777, 2)
```

✓ Check duplicated



```
1 df['English'].value_counts()
```



English	count
i can't tell you how happy i am that you've come to visit us.	29
you're the teacher.	26
you gotta get more organized.	26
get away!	25
you're the oldest.	24
...	...
i know that i didn't do that.	1
i know someone you work with.	1

```
[ ] 1 df_unique = df[  
2     ~df.duplicated(subset=["English"], keep=False) & # Remove if English is duplicated  
3     ~df.duplicated(subset=["French"], keep=False)    # Remove if French is duplicated  
4 ]  
5  
6 print(f"Dataset size after removing all duplicates from both columns: {len(df_unique)}")
```

→ Dataset size after removing all duplicates from both columns: 96873

```
[ ] 1 df_unique.shape
```

→ (96873, 2)

93 exhale. expire.

```
[ ] 1 print(f"Remaining duplicate English values: {df_unique['English'].duplicated().sum()}")
2 print(f"Remaining duplicate French values: {df_unique['French'].duplicated().sum()}")
```

→ Remaining duplicate English values: 0
Remaining duplicate French values: 0

```
[ ] 1 df_unique['English'].value_counts()
```

→

English	count
i went drinking with one of my boyfriend's friends, and now he's furious at me. "was this friend a guy or a girl?" "a guy, obviously. why would i go drinking with his female friends?" "yeah, you're right." "his name is tom. he's really hot, and i really want to go drinking with him again."	1
who?	1
fire!	1
i won.	1

- ▼ Remove Special Characters Only For Better Performance

```
[ ] 1 import re
2
3 # Function to remove special characters
4 def clean_text(text):
5     if isinstance(text, str): # Ensure it's a string before processing
6         return re.sub(r"[^a-zA-ZÀ-ÖØ-Ӯ\s']", "", text) # Keep only letters, spaces, and apostrophes
7     return text
8
9 # Apply function to both columns
10 df_unique["English"] = df_unique["English"].apply(clean_text)
11 df_unique["French"] = df_unique["French"].apply(clean_text)
12
```

Check Similarity of Each Rows English to French are Correctly translated or Correct

Meanings are there

Use Cosine Similarity Between Sentence Embeddings

Check how similar the English and French sentences are using sentence embeddings (from sentence-transformers).

```
[ ]    1 from sentence_transformers import SentenceTransformer  
  2 from sklearn.metrics.pairwise import cosine_similarity  
  3
```

→ Number of low-quality translations: 150

```
[ ] 1 # Keep only high-quality translations (similarity score >= 0.5)
2 df_cleaned = df_sampled[similarity_scores >= 0.5]
3
4 print(f"Dataset after removing low-quality translations: {len(df_cleaned)} rows.")
5
6 # Save the cleaned dataset
7 df_cleaned.to_csv("cleaned_translations.csv", index=False, encoding="utf-8")
8 print("Cleaned dataset saved as 'cleaned_translations.csv'")
```

→ Dataset after removing low-quality translations: 9850 rows.

Cleaned dataset saved as 'cleaned_translations.csv'

▼ BLEU Score is 0.0000

Afte BPE Check BLEW Score

```
[ ] 1 #####
```

▼ Tokenization & Embedding

Feature	BPE	WordPiece
Used In	GPT-2, MarianMT, M2M-100	BERT, mBART, T5
Splitting Strategy	Merges frequent character pairs	Probabilistically splits words
Efficiency	Works well for morphologically rich languages (French, Tamil)	Better for context-sensitive embeddings
Vocabulary Size	Smaller	Larger (handles OOV words better)
Computational Cost	Faster	Slightly slower



English

French

English_Tokens

French_Tokens

English_BPE

French_BPE

0	the mayor provided me with an identity card	le maire m'a fourni une carte d'identité	['the', 'mayor', 'provided', 'me', 'with', 'an...']	['le', 'maire', 'm\'a', 'fourni', 'une', 'carte...']	_the _mayor _provided _me _with _an _identity ...	_le _maire _m' a _fourni _une _carte _d' ide...
1	i called him tom	je l'ai appelé tom	['i', 'called', 'him', 'tom']	['je', "l'ai", 'appelé', 'tom']	_i _called _him _tom	_je _l' ai _appelé _tom
2	i don't meet him so often	je ne le rencontre pas si souvent	['i', "don't", 'meet', 'him', 'so', 'often']	['je', 'ne', 'le', 'rencontre', 'pas', 'si', '...']	_i _don' t _meet _him _so _often	_je _ne _le _rencontre _pas _si _souvent
3	you'll have to come with me	vous allez devoir venir avec moi	["you'll", 'have', 'to', 'come', 'with', 'me']	['vous', 'allez', 'devoir', 'venir', 'avec', '...']	_you 'll _have _to _come _with _me	_vous _allez _devoir _venir _avec _moi
4	whether shakespeare wrote this poem or not wil...	si shakespeare a écrit ce poème ou pas restera...	['whether', 'shakespeare', 'wrote', 'this', 'p...']	['si', 'shakespeare', 'a', 'écrit', 'ce', 'poè...']	_whether _shakespeare _wrote _this _poem _or _...	_si _shakespeare _a _écrit _ce _poème _ou _pas...

▼ BLEU Score Little Bit Increased :- 0.0146

🚀 Interpretation of BLEU Score

- BLEU = 0 → Bad translation (poor match)
- BLEU < 10 → Very low-quality translation
- BLEU 10-30 → Understandable but not fluent
- BLEU 30-50 → Good translation
- BLEU 50+ → Excellent translation

Ram Issues so passing only 1k instead of 98k data that only accuracy & fluency very low

▼ Next Step: Fine Tuning

- **mBART (Multilingual BART)** – Strong for multilingual translation
- **mT5 (Multilingual T5)** – Good for text generation and translation
- **M2M-100 (Facebook's Many-to-Many model)** – Directly trained for translation

mBART (Multilingual BART)

```
wandb: No netrc file found, creating one.
wandb: Appending key for api.wandb.ai to your netrc file: /root/.netrc
wandb: Currently logged in as: antonychackotc (antonychackotc-guvi-geek-networks) to https://api.wandb.ai. Use `wandb login --relogin` to
Tracking run with wandb version 0.19.8
Run data is saved locally in /content/wandb/run-20250324_072159-ty85qsd7
Syncing run /mbart-finetuned to Weights & Biases \(docs\)
View project at https://wandb.ai/antonychackotc-guvi-geek-networks/huggingface
View run at https://wandb.ai/antonychackotc-guvi-geek-networks/huggingface/runs/ty85qsd7
```

Session Crashed

2) Now go to mT5

- Input: i called him tom
✓ Reference: je l'ai appelé tom
✗ Model Output: <extra_id_0>
- Input: i don't meet him so often
✓ Reference: je ne le rencontre pas si souvent
✗ Model Output: <extra_id_0>
- Input: you'll have to come with me
✓ Reference: vous allez devoir venir avec moi
✗ Model Output: <extra_id_0>

Output Shows <extra_id_0>

**BLEU Score Not increased: Because Fine Tuning
not get enough data to train get accuracy &
Fluency if give enough Data Session will
Crashed on mBART**

So Now Run Pretrained Data

Helsinki-NLP/opus-mt-en-fr

Prediction Run Successful

```
14 # Test translation on new input
15 text = "The weather is nice today, and I want to go for a walk."
16 print("✓ Translation:", translate(text))
17
```

```
↳ /usr/local/lib/python3.11/dist-packages/transformers/models/marian/tokenization_marian.py:175: Us
   warnings.warn("Recommended: pip install sacremoses.")
✓ Translation: Le temps est beau aujourd'hui, et je veux aller faire une promenade.
```

Conclusion:

Hands-On Experience Gained:

- **Fine-Tuning Transformer Models:** Successfully fine-tuned the mBART model for insurance fraud detection, learning real-world NLP model training.
- **Tokenization & Preprocessing:** Used MBartTokenizer and prepared text inputs in the right format, understanding how models process text.
- **Model Inference:** Generated predictions with the fine-tuned model, applying techniques like beam search and early stopping for better outputs.
- **Working in Colab:** Gained practical experience training large models efficiently using Google Colab.

Thank You