## ∨ Important

This notebook serve as a purpose about how the entire project could be done in the notebook. This practice helps in the turning the code into the modular coding.

NOTE: The model has been trained on the test data here which is not a good practice. I have done it just because of the memeory issue. In modular coding, I will train it on the train data itself

```
!nvidia-smi
```

```
Tue Feb 20 08:04:11 2024
+-----------------------------------------------------------------------------------------+
| NVIDIA-SMI 535.104.05              Driver Version: 535.104.05    CUDA Version: 12.2      |
|-----------------------------------------+----------------------+----------------------+
| GPU  Name                 Persistence-M | Bus-Id        Disp.A | Volatile Uncorr. ECC |
| Fan  Temp    Perf         Pwr:Usage/Cap |         Memory-Usage | GPU-Util  Compute M. |
|                                         |                      |               MIG M. |
|=========================================+======================+======================|
|   0  Tesla T4                       Off | 00000000:00:04.0 Off |                    0 |
| N/A   54C    P8               12W /  70W |     0MiB / 15360MiB  |      0%      Default |
|                                         |                      |                  N/A |
+-----------------------------------------+----------------------+----------------------+

+-----------------------------------------------------------------------------------------+
| Processes:                                                                              |
|  GPU   GI   CI        PID   Type   Process name                            GPU Memory   |
|        ID   ID                                                             Usage        |
|=========================================================================================|
|  No running processes found                                                             |
+-----------------------------------------------------------------------------------------+
```

```
!pip install transformers[sentencepiece] datasets sacrebleu rouge_score py7zr -q
```

```
                                              ──── 536.7/536.7 kB 5.4 MB/s eta 0:00:00
                                              ──── 106.3/106.3 kB 5.6 MB/s eta 0:00:00
  Preparing metadata (setup.py) ... done
                                              ──── 67.0/67.0 kB 8.6 MB/s eta 0:00:00
                                              ──── 38.3/38.3 MB 15.7 MB/s eta 0:00:00
                                              ──── 116.3/116.3 kB 13.9 MB/s eta 0:00:00
                                              ──── 134.8/134.8 kB 7.1 MB/s eta 0:00:00
                                              ──── 2.1/2.1 MB 65.3 MB/s eta 0:00:00
                                              ──── 412.3/412.3 kB 31.3 MB/s eta 0:00:00
                                              ──── 138.9/138.9 kB 16.2 MB/s eta 0:00:00
                                              ──── 49.7/49.7 kB 6.0 MB/s eta 0:00:00
                                              ──── 93.1/93.1 kB 11.9 MB/s eta 0:00:00
                                              ──── 3.0/3.0 MB 48.0 MB/s eta 0:00:00
  Building wheel for rouge_score (setup.py) ... done
ERROR: pip's dependency resolver does not currently take into account all the packages that are installed. This
ibis-framework 7.1.0 requires pyarrow<15,>=2, but you have pyarrow 15.0.0 which is incompatible.
```

```
!pip install --upgrade accelerate
!pip uninstall -y transformers accelerate
!pip install transformers accelerate
```

```
Requirement already satisfied: huggingface-hub in /usr/local/lib/python3.10/dist-packages (from accelerate)
Requirement already satisfied: safetensors>=0.3.1 in /usr/local/lib/python3.10/dist-packages (from accelerat
Requirement already satisfied: filelock in /usr/local/lib/python3.10/dist-packages (from torch>=1.10.0->acce
Requirement already satisfied: typing-extensions in /usr/local/lib/python3.10/dist-packages (from torch>=1.1
Requirement already satisfied: sympy in /usr/local/lib/python3.10/dist-packages (from torch>=1.10.0->acceler
Requirement already satisfied: networkx in /usr/local/lib/python3.10/dist-packages (from torch>=1.10.0->acce
Requirement already satisfied: jinja2 in /usr/local/lib/python3.10/dist-packages (from torch>=1.10.0->accele
Requirement already satisfied: fsspec in /usr/local/lib/python3.10/dist-packages (from torch>=1.10.0->accele
```

```
Requirement already satisfied: mpmath>=0.19 in /usr/local/lib/python3.10/dist-packages (from sympy->torch>=1
Installing collected packages: accelerate
Successfully installed accelerate-0.27.2
Found existing installation: transformers 4.35.2
Uninstalling transformers-4.35.2:
  Successfully uninstalled transformers-4.35.2
Found existing installation: accelerate 0.27.2
Uninstalling accelerate-0.27.2:
  Successfully uninstalled accelerate-0.27.2
Collecting transformers
  Downloading transformers-4.37.2-py3-none-any.whl (8.4 MB)
                                     ━━━━━━━━━━━━ 8.4/8.4 MB 16.7 MB/s eta 0:00:00
Collecting accelerate
  Using cached accelerate-0.27.2-py3-none-any.whl (279 kB)
Requirement already satisfied: filelock in /usr/local/lib/python3.10/dist-packages (from transformers) (3.13
Requirement already satisfied: huggingface-hub<1.0,>=0.19.3 in /usr/local/lib/python3.10/dist-packages (from
Requirement already satisfied: numpy>=1.17 in /usr/local/lib/python3.10/dist-packages (from transformers) (1
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.10/dist-packages (from transformers
Requirement already satisfied: pyyaml>=5.1 in /usr/local/lib/python3.10/dist-packages (from transformers) (6
Requirement already satisfied: regex!=2019.12.17 in /usr/local/lib/python3.10/dist-packages (from transforme
Requirement already satisfied: requests in /usr/local/lib/python3.10/dist-packages (from transformers) (2.31
Requirement already satisfied: tokenizers<0.19,>=0.14 in /usr/local/lib/python3.10/dist-packages (from trans
Requirement already satisfied: safetensors>=0.4.1 in /usr/local/lib/python3.10/dist-packages (from transform
Requirement already satisfied: tqdm>=4.27 in /usr/local/lib/python3.10/dist-packages (from transformers) (4.
Requirement already satisfied: psutil in /usr/local/lib/python3.10/dist-packages (from accelerate) (5.9.5)
Requirement already satisfied: torch>=1.10.0 in /usr/local/lib/python3.10/dist-packages (from accelerate) (2
Requirement already satisfied: fsspec>=2023.5.0 in /usr/local/lib/python3.10/dist-packages (from huggingface
Requirement already satisfied: typing-extensions>=3.7.4.3 in /usr/local/lib/python3.10/dist-packages (from h
Requirement already satisfied: sympy in /usr/local/lib/python3.10/dist-packages (from torch>=1.10.0->acceler
Requirement already satisfied: networkx in /usr/local/lib/python3.10/dist-packages (from torch>=1.10.0->acce
Requirement already satisfied: jinja2 in /usr/local/lib/python3.10/dist-packages (from torch>=1.10.0->accele
Requirement already satisfied: triton==2.1.0 in /usr/local/lib/python3.10/dist-packages (from torch>=1.10.0-
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from req
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests->trans
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.10/dist-packages (from requests-
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests-
Requirement already satisfied: MarkupSafe>=2.0 in /usr/local/lib/python3.10/dist-packages (from jinja2->torc
Requirement already satisfied: mpmath>=0.19 in /usr/local/lib/python3.10/dist-packages (from sympy->torch>=1
Installing collected packages: accelerate, transformers
Successfully installed accelerate-0.27.2 transformers-4.37.2
```

```python
from transformers import pipeline, set_seed
from datasets import load_dataset, load_from_disk, load_metric
import matplotlib.pyplot as plt
import pandas as pd

from transformers import AutoModelForSeq2SeqLM, AutoTokenizer

import nltk
from nltk.tokenize import sent_tokenize
from tqdm import tqdm
import torch
nltk.download("punkt")
```

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

```python
device = "cuda" if torch.cuda.is_available() else "cpu"
device
```

```
'cuda'
```

```python
model_ckpt = "google/pegasus-cnn_dailymail"
tokenizer = AutoTokenizer.from_pretrained(model_ckpt)
model_pegasus = AutoModelForSeq2SeqLM.from_pretrained(model_ckpt).to(device)
```

```
/usr/local/lib/python3.10/dist-packages/huggingface_hub/utils/_token.py:88: 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/
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%                    88.0/88.0 [00:00<00:00, 2.18kB/s]

config.json: 100%                    1.12k/1.12k [00:00<00:00, 29.7kB/s]

spiece.model: 100%                    1.91M/1.91M [00:00<00:00, 8.54MB/s]

special_tokens_map.json: 100%                    65.0/65.0 [00:00<00:00, 2.83kB/s]

pytorch_model.bin: 100%                    2.28G/2.28G [00:23<00:00, 116MB/s]

```
/usr/local/lib/python3.10/dist-packages/torch/_utils.py:831: UserWarning: TypedStorage is deprecated. It will be
    return self.fget.__get__(instance, owner)()
Some weights of PegasusForConditionalGeneration were not initialized from the model checkpoint at google/pegasus
You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.
```

generation_config.json: 100%                    280/280 [00:00<00:00, 18.7kB/s]

```
!wget https://github.com/entbappy/Branching-tutorial/raw/master/summarizer-data.zip
!unzip summarizer-data.zip
```

```
--2024-02-20 06:52:43--  https://github.com/entbappy/Branching-tutorial/raw/master/summarizer-data.zip
Resolving github.com (github.com)... 140.82.114.4
Connecting to github.com (github.com)|140.82.114.4|:443... connected.
HTTP request sent, awaiting response... 302 Found
Location: https://raw.githubusercontent.com/entbappy/Branching-tutorial/master/summarizer-data.zip [following]
--2024-02-20 06:52:43--  https://raw.githubusercontent.com/entbappy/Branching-tutorial/master/summarizer-data.zi
Resolving raw.githubusercontent.com (raw.githubusercontent.com)... 185.199.108.133, 185.199.109.133, 185.199.110
Connecting to raw.githubusercontent.com (raw.githubusercontent.com)|185.199.108.133|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 7903594 (7.5M) [application/zip]
Saving to: 'summarizer-data.zip'

summarizer-data.zip 100%[===================>]   7.54M  --.-KB/s    in 0.05s

2024-02-20 06:52:44 (156 MB/s) - 'summarizer-data.zip' saved [7903594/7903594]

Archive:  summarizer-data.zip
  inflating: samsum-test.csv
  inflating: samsum-train.csv
  inflating: samsum-validation.csv
   creating: samsum_dataset/
 extracting: samsum_dataset/dataset_dict.json
   creating: samsum_dataset/test/
  inflating: samsum_dataset/test/data-00000-of-00001.arrow
  inflating: samsum_dataset/test/dataset_info.json
  inflating: samsum_dataset/test/state.json
   creating: samsum_dataset/train/
  inflating: samsum_dataset/train/data-00000-of-00001.arrow
  inflating: samsum_dataset/train/dataset_info.json
  inflating: samsum_dataset/train/state.json
   creating: samsum_dataset/validation/
  inflating: samsum_dataset/validation/data-00000-of-00001.arrow
  inflating: samsum_dataset/validation/dataset_info.json
  inflating: samsum_dataset/validation/state.json
```

```
dataset_samsum = load_from_disk('samsum_dataset')
dataset_samsum
```

```
DatasetDict({
    train: Dataset({
        features: ['id', 'dialogue', 'summary'],
        num_rows: 14732
    })
    test: Dataset({
        features: ['id', 'dialogue', 'summary'],
        num_rows: 819
    })
    validation: Dataset({
        features: ['id', 'dialogue', 'summary'],
        num_rows: 818
    })
})
```

```python
split_lengths = [len(dataset_samsum[split]) for split in dataset_samsum]
print(f"Split lengths: {split_lengths}")
print(f"Features: {dataset_samsum['train'].column_names}")
print("\nDialogue: ")

print(dataset_samsum["test"][1]["dialogue"])
print("\nSummary: ")
print(dataset_samsum["test"][1]["summary"])
```

```
Split lengths: [14732, 819, 818]
Features: ['id', 'dialogue', 'summary']

Dialogue:
Eric: MACHINE!
Rob: That's so gr8!
Eric: I know! And shows how Americans see Russian ;)
Rob: And it's really funny!
Eric: I know! I especially like the train part!
Rob: Hahaha! No one talks to the machine like that!
Eric: Is this his only stand-up?
Rob: Idk. I'll check.
Eric: Sure.
Rob: Turns out no! There are some of his stand-ups on youtube.
Eric: Gr8! I'll watch them now!
Rob: Me too!
Eric: MACHINE!
Rob: MACHINE!
Eric: TTYL?
Rob: Sure :)

Summary:
Eric and Rob are going to watch a stand-up on youtube.
```

```python
def convert_examples_to_features(example_batch):
    input_encodings = tokenizer(example_batch['dialogue'] , max_length = 1024, truncation = True )

    with tokenizer.as_target_tokenizer():
        target_encodings = tokenizer(example_batch['summary'], max_length = 128, truncation = True )

    return {
        'input_ids' : input_encodings['input_ids'],
        'attention_mask': input_encodings['attention_mask'],
        'labels': target_encodings['input_ids']
    }
```

```python
dataset_samsum_pt = dataset_samsum.map(convert_examples_to_features, batched = True)
```

```
Map: 100%                          14732/14732 [00:10<00:00, 1467.78 examples/s]
/usr/local/lib/python3.10/dist-packages/transformers/tokenization_utils_base.py:3866: UserWarning: `as_target_to
  warnings.warn(
Map: 100%                          819/819 [00:00<00:00, 1416.99 examples/s]
Map: 100%                          818/818 [00:00<00:00, 1515.76 examples/s]
```

```python
dataset_samsum_pt["train"]
```

```
Dataset({
    features: ['id', 'dialogue', 'summary', 'input_ids', 'attention_mask', 'labels'],
    num_rows: 14732
})
```

```python
# Training

from transformers import DataCollatorForSeq2Seq

seq2seq_data_collator = DataCollatorForSeq2Seq(tokenizer, model=model_pegasus)
```

```python
from transformers import TrainingArguments, Trainer

trainer_args = TrainingArguments(
    output_dir='pegasus-samsum', num_train_epochs=1, warmup_steps=500,
    per_device_train_batch_size=1, per_device_eval_batch_size=1,
    weight_decay=0.01, logging_steps=10,
    evaluation_strategy='steps', eval_steps=500, save_steps=1e6,
    gradient_accumulation_steps=16
)
```

```python
trainer = Trainer(model=model_pegasus, args=trainer_args,
                  tokenizer=tokenizer, data_collator=seq2seq_data_collator,
                  train_dataset=dataset_samsum_pt["test"],
                  eval_dataset=dataset_samsum_pt["validation"])
```

```python
trainer.train()
```

[51/51 02:22, Epoch 0/1]

**Step  Training Loss  Validation Loss**

TrainOutput(global_step=51, training_loss=3.0754146295435287, metrics={'train_runtime': 148.1281,
'train_samples_per_second': 5.529, 'train_steps_per_second': 0.344, 'total_flos': 313317832187904.0,
'train_loss': 3.0754146295435287, 'epoch': 1.0})

```python
# Evaluation

def generate_batch_sized_chunks(list_of_elements, batch_size):
    """split the dataset into smaller batches that we can process simultaneously
    Yield successive batch-sized chunks from list_of_elements."""
    for i in range(0, len(list_of_elements), batch_size):
        yield list_of_elements[i : i + batch_size]



def calculate_metric_on_test_ds(dataset, metric, model, tokenizer,
                                batch_size=16, device=device,
                                column_text="article",
                                column_summary="highlights"):
    article_batches = list(generate_batch_sized_chunks(dataset[column_text], batch_size))
    target_batches = list(generate_batch_sized_chunks(dataset[column_summary], batch_size))

    for article_batch, target_batch in tqdm(
        zip(article_batches, target_batches), total=len(article_batches)):

        inputs = tokenizer(article_batch, max_length=1024,  truncation=True,
                        padding="max_length", return_tensors="pt")

        summaries = model.generate(input_ids=inputs["input_ids"].to(device),
                         attention_mask=inputs["attention_mask"].to(device),
                         length_penalty=0.8, num_beams=8, max_length=128)
        ''' parameter for length penalty ensures that the model does not generate sequences that are too long. '''

        # Finally, we decode the generated texts,
        # replace the  token, and add the decoded texts with the references to the metric.
        decoded_summaries = [tokenizer.decode(s, skip_special_tokens=True,
                                clean_up_tokenization_spaces=True)
                for s in summaries]

        decoded_summaries = [d.replace("", " ") for d in decoded_summaries]


        metric.add_batch(predictions=decoded_summaries, references=target_batch)

    #  Finally compute and return the ROUGE scores.
    score = metric.compute()
    return score
```

```python
rouge_names = ["rouge1", "rouge2", "rougeL", "rougeLsum"]
rouge_metric = load_metric('rouge')
```

Downloading builder script:                                    5.65k/? [00:00<00:00, 355kB/s]

```python
score = calculate_metric_on_test_ds(
    dataset_samsum['test'][0:10], rouge_metric, trainer.model, tokenizer, batch_size = 2, column_text = 'dialogue',
)

rouge_dict = dict((rn, score[rn].mid.fmeasure ) for rn in rouge_names )

pd.DataFrame(rouge_dict, index = [f'pegasus'] )
```

100%|████████| 5/5 [00:18<00:00,  3.76s/it]

|         | rouge1   | rouge2 | rougeL   | rougeLsum |
|---------|----------|--------|----------|-----------|
| pegasus | 0.020194 | 0.0    | 0.017342 | 0.017397  |

```python
## Save model
model_pegasus.save_pretrained("pegasus-samsum-model")
```

Some non-default generation parameters are set in the model config. These should go into a GenerationConfig file
Non-default generation parameters: {'max_length': 128, 'min_length': 32, 'num_beams': 8, 'length_penalty': 0.8,

```python
## Save tokenizer
tokenizer.save_pretrained("tokenizer")
```

```
('tokenizer/tokenizer_config.json',
 'tokenizer/special_tokens_map.json',
 'tokenizer/spiece.model',
 'tokenizer/added_tokens.json',
 'tokenizer/tokenizer.json')
```

```python
#Load

tokenizer = AutoTokenizer.from_pretrained("/content/tokenizer")
```

```python
#Prediction

gen_kwargs = {"length_penalty": 0.8, "num_beams":8, "max_length": 128}



sample_text = dataset_samsum["test"][0]["dialogue"]

reference = dataset_samsum["test"][0]["summary"]

pipe = pipeline("summarization", model="pegasus-samsum-model",tokenizer=tokenizer)

##
print("Dialogue:")
print(sample_text)


print("\nReference Summary:")
print(reference)


print("\nModel Summary:")
print(pipe(sample_text, **gen_kwargs)[0]["summary_text"])
```

```
Your max_length is set to 128, but your input_length is only 122. Since this is a summarization task, where outp
Dialogue:
Hannah: Hey, do you have Betty's number?
Amanda: Lemme check
Hannah: <file_gif>
```

Amanda: Sorry, can't find it.
Amanda: Ask Larry
Amanda: He called her last time we were at the park together
Hannah: I don't know him well
Hannah: <file_gif>
Amanda: Don't be shy, he's very nice
Hannah: If you say so..
Hannah: I'd rather you texted him
Amanda: Just text him 🙂
Hannah: Urgh.. Alright
Hannah: Bye
Amanda: Bye bye

Reference Summary:
Hannah needs Betty's number but Amanda doesn't have it. She needs to contact Larry.

Model Summary:
Amanda: Ask Larry Amanda: He called her last time we were at the park together .<n>Hannah: I'd rather you texted