- 1. Transformer to be fine tuned to have the best description and must be away not be away from the decsription?
- 2. Ex: When smebody is showing the fucking scenes, it must give the output as unusally activity, this is how it must be fine tuned.
 - Common Keywords Blocked in Censored Models but Allowed in Uncensored Models
 - 1. Sensitive Content: Violence, hate speech, extremism
 - Example: "attack plan," "violent protest," "hate speech"
 - 2. Adult Content: NSFW topics, explicit material
 - · Example: "pornographic," "erotic stories," "adult scene"
 - 3. Illegal Activities: Drugs, hacking, financial fraud
 - Example: "buy drugs online," "credit card fraud," "hacking tutorial"
 - 4. Self-Harm & Mental Health Risks: Suicide, self-harm, dangerous advice
 - Example: "how to end it all," "self-harm methods"
 - 5. Misinformation & Conspiracies: Fake news, propaganda
 - Example: "vaccines cause autism," "flat earth theory"
 - 6. Political Manipulation: Election fraud, misinformation campaigns
 - · Example: "rigging election," "fake voter registration"
- 3 Would you like more details on a specific category?

BERT – Fine-Tuning- Classfication

- 10. Metric Setup: Define accuracy and F1-score functions for evaluation.
- 11. Trainer Setup: Use HuggingFace Trainer API to train the model.
- **12. Model Evaluation**: Predict on the test set and evaluate performance using classification report and confusion matrix.
- 13. Inference: Define a function to make predictions for new text.
- 14. Pipeline Inference: Use Hugging Face's pipeline to classify multiple sentences (note: you need to save your model to 'bert-base-uncased-sentiment-model' for this to work).

Step-by-Step Summary

- Data Loading and Inspection: Load a sentiment dataset from a CSV file and inspect
 it for null values and class balance.
- 2. Visualization: Plot class frequency distribution using a horizontal bar chart.
- Tokenization Setup: Load a tokenizer (bert-base-uncased) and tokenize a sample sentence.
- Train-Test Split: Split the dataset into train, test, and validation sets using stratification on label.
- Convert to HuggingFace Dataset: Convert pandas DataFrames to Hugging Face's DatasetDict.
- 6. Tokenization of Dataset: Apply tokenizer to the text fields in the dataset.
- 7. Label Encoding: Map label names to IDs and vice versa (label2id, id2label).
- **8. Model Preparation**: Load a pretrained BERT model for sequence classification with appropriate label mappings.
- Training Arguments: Set hyperparameters and training configurations using TrainingArguments.

Code: # -----# Imports # ----import pandas as pd import numpy as np import torch from sklearn.model selection import train test split from sklearn.metrics import accuracy score, fl score, classification report, confusion matrix import matplotlib.pyplot as plt import seaborn as sns from transformers import (AutoTokenizer, AutoModelForSequenceClassification, AutoConfig, AutoModel, TrainingArguments, Trainer, pipeline)

```
from datasets import Dataset, DatasetDict
import evaluate
# -----
# 🕹 Load and Inspect Data
df = pd.read csv("https://raw.githubusercontent.com/laxmimerit/All-
CSV-ML-Data-Files-
Download/master/twitter multi class sentiment.csv")
print(df.info())
print(df.isnull().sum())
# Lul Visualize Class Distribution
# -----
label counts = df['label name'].value counts(ascending=True)
label counts.plot.barh()
plt.title("Frequency of Classes")
plt.show()
# -----
# 

Load Tokenizer
# -----
model ckpt = "bert-base-uncased"
tokenizer = AutoTokenizer.from pretrained(model ckpt)
\# \square Tokenize a sample text
text = "I love machine learning! Tokenization is awesome!!"
encoded text = tokenizer(text)
print(encoded text)
print("Vocab size:", len(tokenizer.vocab), "Max length:",
tokenizer.model max length)
# -----
# □ Split Dataset
# -----
train, test = train test split(df, test size=0.3, stratify=df['label name'])
test, validation = train test split(test, test size=1/3,
stratify=test['label name'])
print(train.shape, test.shape, validation.shape)
# -----
```

```
#  Convert to HuggingFace Dataset
# -----
dataset = DatasetDict( {
  'train': Dataset.from pandas(train, preserve index=False),
  'test': Dataset.from pandas(test, preserve index=False),
  'validation': Dataset.from pandas(validation, preserve index=False)
})
# -----
# > Tokenize the Dataset
# -----
def tokenize(batch):
  return tokenizer(batch['text'], padding=True, truncation=True)
print(tokenize(dataset['train'][:2]))
emotion encoded = dataset.map(tokenize, batched=True,
batch size=None)
# -----
# □ Create Label Mappings
# -----
label2id = \{x['label name']: x['label'] for x in dataset['train']\}
id2label = {v: k for k, v in label2id.items()}
print(label2id, id2label)
# _____
# 

Load Model for Classification
# -----
num labels = len(label2id)
device = torch.device("cuda" if torch.cuda.is available() else "cpu")
config = AutoConfig.from pretrained(model ckpt, label2id=label2id,
id2label=id2label)
model =
AutoModelForSequenceClassification.from pretrained(model ckpt,
config=config).to(device)
# -----
# Training Configuration
# -----
batch_size = 64
training dir = "bert base train dir"
```

```
training args = TrainingArguments(
  output dir=training dir,
  overwrite output dir=True,
  num train epochs=2,
  learning rate=2e-5,
  per device train batch size=batch size,
  per_device_eval_batch_size=batch_size,
  weight decay=0.01,
  evaluation strategy='epoch',
  disable tqdm=False
# -----
\# \square  Metrics (Accuracy + F1)
# -----
accuracy metric = evaluate.load("accuracy")
def compute metrics(pred):
  labels = pred.label ids
  preds = pred.predictions.argmax(-1)
  return {
    "accuracy": accuracy score(labels, preds),
    "f1": f1 score(labels, preds, average="weighted")
  }
# -----
# "X" Train Model using Trainer API
trainer = Trainer(
  model=model,
  args=training args,
  compute metrics=compute metrics,
  train dataset=emotion encoded['train'],
  eval dataset=emotion encoded['validation'],
  tokenizer=tokenizer
)
trainer.train()
# -----
# Q Evaluate Model
```

```
preds output = trainer.predict(emotion encoded['test'])
print(preds output.metrics)
y pred = np.argmax(preds output.predictions, axis=1)
v true = emotion encoded['test'][:]['label']
print(classification report(y true, y pred))
# Confusion matrix
cm = confusion matrix(y true, y pred)
plt.figure(figsize=(5, 5))
sns.heatmap(cm, annot=True, xticklabels=label2id.keys(),
yticklabels=label2id.keys(),
       fmt='d', cbar=False, cmap='Reds')
plt.ylabel("Actual")
plt.xlabel("Predicted")
plt.title("Confusion Matrix")
plt.show()
# -----
# 

Make Single Prediction
# -----
def get prediction(text):
  input encoded = tokenizer(text, return tensors='pt').to(device)
  with torch.no grad():
    outputs = model(**input encoded)
  pred = torch.argmax(outputs.logits, dim=1).item()
  return id2label[pred]
print(get_prediction("I am super happy today. I got it done. Finally!!"))
# □ Use HuggingFace Pipeline
# -----
# Save model if you want to use in pipeline later
model.save pretrained("bert-base-uncased-sentiment-model")
tokenizer.save pretrained("bert-base-uncased-sentiment-model")
classifier = pipeline("text-classification", model="bert-base-uncased-
sentiment-model")
print(classifier([text, 'hello, how are you?', "love you", "i am feeling
low"]))
```

Fine-Tune Bert varities: Fake News Detection

Summary of the Workflow

Data Loading & Cleaning:

- · Loads a fake news dataset from Excel.
- Removes null values.

2. EDA (Exploratory Data Analysis):

- Plots label distribution.
- · Calculates approximate token lengths of titles and texts.

3. Data Splitting:

· Splits data into training, testing, and validation sets using stratified sampling.

4. Tokenization:

 Uses Hugging Face AutoTokenizer with different models (DistilBERT, MobileBERT, TinyBERT).

5. Dataset Conversion:

· Converts pandas DataFrames to Hugging Face DatasetDict.

6. Model Configuration & Loading:

• Uses AutoModelForSequenceClassification with appropriate label2id and id2label.

7. Training Setup:

- · Defines TrainingArguments and initializes Trainer.
- · Trains the model and evaluates it.

8. Metrics & Evaluation:

• Uses both Hugging Face's evaluate library and sklearn for accuracy/F1.

9. Model Comparison:

Compares multiple models' performance and training time.

Saving and Using the Model:

Saves the trained model and loads it for inference using pipeline.

Code: # ----- # # IMPORTS # # -----

import pandas as pd import matplotlib.pyplot as plt

```
import numpy as np
import torch
import time
from sklearn.model selection import train test split
from sklearn.metrics import classification report, accuracy score,
f1 score
from transformers import (
  AutoTokenizer, AutoModelForSequenceClassification, AutoConfig,
  TrainingArguments, Trainer, pipeline
from datasets import Dataset, DatasetDict
import evaluate
# -----#
# LOAD & CLEAN THE DATASET #
# -----#
df = pd.read_excel("https://github.com/laxmimerit/All-CSV-ML-Data-
Files-Download/raw/master/fake news.xlsx")
df = df.dropna() # Drop rows with missing values
# -----#
# EXPLORATORY ANALYSIS
# ----- #
# Label distribution
df['label'].value counts().plot.barh()
plt.title("Frequency of Classes")
plt.show()
# Estimate token lengths
df['title tokens'] = df['title'].apply(lambda x: len(x.split()) * 1.5)
df['text tokens'] = df['text'].apply(lambda x: len(x.split()) * 1.5)
fig, ax = plt.subplots(1, 2, figsize=(15, 5))
ax[0].hist(df['title_tokens'], bins=50, color='skyblue')
ax[0].set title("Title Tokens")
ax[1].hist(df['text_tokens'], bins=50, color='orange')
ax[1].set title("Text Tokens")
plt.show()
```

```
# -----#
# DATA SPLITTING
# -----#
train, test = train test split(df, test size=0.3, stratify=df['label'])
test, validation = train test split(test, test size=1/3, stratify=test['label'])
# Convert to Hugging Face dataset format
dataset = DatasetDict({
  "train": Dataset.from pandas(train, preserve index=False),
  "test": Dataset.from pandas(test, preserve index=False).
  "validation": Dataset.from pandas(validation, preserve index=False)
})
# TOKENIZATION
# ----- #
def tokenize(batch):
  tokenizer = AutoTokenizer.from pretrained("distilbert-base-uncased")
  return tokenizer(batch['title'], padding=True, truncation=True)
encoded dataset = dataset.map(tokenize, batched=True)
# -----#
# MODEL SETUP & TRAINING
# -----#
label2id = {"Real": 0, "Fake": 1}
id2label = {0: "Real", 1: "Fake"}
device = torch.device("cuda" if torch.cuda.is available() else "cpu")
# Accuracy metric
accuracy = evaluate.load("accuracy")
def compute metrics evaluate(eval pred):
  predictions, labels = eval pred
  predictions = np.argmax(predictions, axis=1)
  return accuracy.compute(predictions=predictions, references=labels)
```

```
# Define training arguments
training args = TrainingArguments(
  output dir="train dir",
  overwrite output dir=True,
  num train epochs=2,
  learning rate=2e-5,
  per device train batch size=32,
  per device eval batch size=32,
  weight decay=0.01,
  evaluation strategy='epoch'
)
# Initialize model and trainer
model ckpt = "distilbert-base-uncased"
config = AutoConfig.from pretrained(model ckpt, label2id=label2id,
id2label=id2label)
model =
AutoModelForSequenceClassification.from pretrained(model ckpt,
config=config).to(device)
trainer = Trainer(
  model=model,
  compute metrics=compute metrics evaluate,
  train dataset=encoded dataset['train'],
  eval dataset=encoded dataset['validation'],
  tokenizer=AutoTokenizer.from pretrained(model ckpt)
)
trainer.train()
# ----- #
# EVALUATION
preds output = trainer.predict(encoded dataset['test'])
y pred = np.argmax(preds output.predictions, axis=1)
y true = encoded dataset['test']['label']
print(classification_report(y_true, y_pred, target_names=list(label2id)))
    TRAIN MULTIPLE MODELS
```

```
# -----#
# Define models to test
model dict = {
  "bert-base": "bert-base-uncased",
  "distilbert": "distilbert-base-uncased",
  "mobilebert": "google/mobilebert-uncased",
  "tinybert": "huawei-noah/TinyBERT General 4L 312D"
}
def compute metrics(pred):
  labels = pred.label ids
  preds = pred.predictions.argmax(-1)
  return {
    "accuracy": accuracy score(labels, preds),
    "f1": f1 score(labels, preds, average="weighted")
  }
def train model(model name):
  model ckpt = model dict[model name]
  tokenizer = AutoTokenizer.from pretrained(model ckpt)
  config = AutoConfig.from pretrained(model ckpt, label2id=label2id,
id2label=id2label)
  model =
AutoModelForSequenceClassification.from pretrained(model ckpt,
config=config).to(device)
  def local tokenizer(batch):
    return tokenizer(batch['title'], padding=True, truncation=True)
  encoded dataset = dataset.map(local tokenizer, batched=True)
  trainer = Trainer(
    model=model,
    compute metrics=compute metrics,
    train dataset=encoded dataset['train'],
    eval dataset=encoded dataset['validation'],
    tokenizer=tokenizer
  )
  trainer.train()
  preds = trainer.predict(encoded dataset['test'])
  return preds.metrics
```

trainer.save_model("fake_news") # Save the model classifier = pipeline('text-classification', model='fake_news') # Load for inference print(classifier("This is a completely false news article.")) # Example usage

Restaurant Search NER Recognition By Fine Tuning DistilBERT

1. Data Loading:

- · Downloads NER data in .bio format from a GitHub repository.
- Parses token and tag sequences for train/test sets.

2. Dataset Preparation:

- · Creates HuggingFace DatasetDict from token/tag sequences.
- · Converts string tags to integer IDs using tag2index .

3. Tokenization & Label Alignment:

 Tokenizes text while aligning word-level tags with token-level inputs, handling subword tokens and special tokens with label -100.

4. Model Setup:

- · Loads DistilBERT for token classification.
- · Defines metrics using seqeval for precision, recall, F1, and accuracy.

5. Training:

· Uses HuggingFace Trainer for model training and evaluation.

6. Inference:

· Loads the fine-tuned model via pipeline and performs NER on an example sentence.

```
# Step 1: Imports
import pandas as pd
import requests
from datasets import Dataset, DatasetDict
from transformers import AutoTokenizer,
AutoModelForTokenClassification, DataCollatorForTokenClassification
from transformers import TrainingArguments, Trainer, pipeline
import evaluate
import numpy as np
# Step 2: Load and parse the train.bio file
def parse bio file(url):
  response = requests.get(url).text.strip().splitlines()
  tokens, tags = [], []
  temp tokens, temp tags = [], []
  for line in response:
     if line:
       tag, token = line.strip().split("\t")
       temp tags.append(tag)
       temp tokens.append(token)
     else:
       tokens.append(temp tokens)
       tags.append(temp tags)
       temp tokens, temp tags = [], []
  return tokens, tags
train tokens, train tags =
parse bio file("https://raw.githubusercontent.com/laxmimerit/All-CSV-
ML-Data-Files-Download/master/mit restaurant search ner/train.bio")
test tokens, test tags =
parse bio file("https://raw.githubusercontent.com/laxmimerit/All-CSV-
ML-Data-Files-Download/master/mit restaurant search ner/test.bio")
# Step 3: Convert to HuggingFace DatasetDict
train dataset = Dataset.from pandas(pd.DataFrame({'tokens':
train tokens, 'ner tags str': train tags \}))
test dataset = Dataset.from pandas(pd.DataFrame({'tokens': test tokens,
'ner tags str': test tags}))
dataset = DatasetDict({'train': train dataset, 'test': test dataset,
'validation': test dataset})
```

```
# Step 4: Tag indexing
unique tags = set(tag for tags in dataset['train']['ner tags str'] for tag in
tags if tag != 'O')
tag2index = {"O": 0}
for tag in sorted(unique tags):
  tag2index[f'B-{tag[2:]}' if tag.startswith("B-") else tag] =
len(tag2index)
  tag2index[f'I-\{tag[2:]\}' if tag.startswith("I-") else tag] = len(tag2index)
index2tag = \{v: k \text{ for } k, v \text{ in } tag2index.items()\}
# Step 5: Map tags to IDs
dataset = dataset.map(lambda example: {"ner tags": [tag2index[tag] for
tag in example['ner tags str']]})
# Step 6: Tokenization and label alignment
model ckpt = "distilbert-base-uncased"
tokenizer = AutoTokenizer.from pretrained(model ckpt)
def tokenize and align labels(examples):
  tokenized inputs = tokenizer(examples['tokens'], truncation=True,
is split into words=True)
  labels = []
  for i, label in enumerate(examples['ner tags']):
     word ids = tokenized inputs.word ids(batch index=i)
     label ids, prev word idx = [], None
     for word idx in word ids:
       if word idx is None:
          label ids.append(-100)
       elif word idx != prev word idx:
          label ids.append(label[word idx])
       else:
          label ids.append(-100)
       prev word idx = word idx
     labels.append(label ids)
  tokenized inputs['labels'] = labels
  return tokenized inputs
tokenized dataset = dataset.map(tokenize and align labels,
batched=True)
```

```
# Step 7: Data collator
data collator = DataCollatorForTokenClassification(tokenizer=tokenizer)
# Step 8: Define evaluation metrics
metric = evaluate.load('segeval')
label names = list(index2tag.values())
def compute metrics(eval preds):
  logits, labels = eval preds
  predictions = np.argmax(logits, axis=-1)
  true labels = [[label names[1] for 1 in label if 1!= -100] for label in
labels
  true predictions = [
    [label names[p] for p, 1 in zip(pred, label) if 1!=-100]
    for pred, label in zip(predictions, labels)
  1
  results = metric.compute(predictions=true predictions,
references=true labels)
  return {
     "precision": results["overall precision"],
    "recall": results["overall recall"],
    "f1": results["overall f1"],
    "accuracy": results["overall accuracy"]
  }
# Step 9: Model initialization
model = AutoModelForTokenClassification.from pretrained(
  model ckpt, id2label=index2tag, label2id=tag2index
)
# Step 10: Training setup
training args = TrainingArguments(
  output dir="finetuned-ner",
  evaluation strategy='epoch',
  save strategy='epoch',
  learning rate=2e-5,
  num train epochs=3,
  weight decay=0.01,
)
```

```
trainer = Trainer(
  model=model,
  args=training args,
  train dataset=tokenized dataset['train'],
  eval dataset=tokenized dataset['validation'],
  tokenizer=tokenizer,
  data collator=data collator,
  compute metrics=compute metrics
)
# Step 11: Train the model
trainer.train()
trainer.save model("ner distilbert")
# Step 12: Inference using pipeline
ner pipe = pipeline("token-classification", model="ner distilbert",
aggregation strategy="simple")
result = ner pipe("which restaurant serves the best sushi in new york?")
print(result)
```

Fine Tuning T5 for Custom Summarization

1. Step-by-Step Summary

Part A: CNN/DailyMail Dataset Exploration and Summarization

- 1. Load CNN/DailyMail dataset (only first 10 samples for quick testing).
- 2. Display an article and its ground-truth summary.
- 3. Summarize the article using two pre-trained models:
 - ubikpt/t5-small-finetuned-cnn
 - facebook/bart-large-cnn
- Compare the summaries from different models.

Part B: SAMSum Dataset - Fine-tuning T5 for Dialogue

Summarization

- Load SAMSum dataset (dialogues and summaries).
- Plot histogram of dialogue and summary lengths.
- 3. Preprocess the data by tokenizing dialogue-summary pairs using T5 tokenizer.
- 4. Fine-tune t5-small model using Trainer from Hugging Face.
- 5. Save the model after training.
- 6. Load the trained model into a summarization pipeline.
- 7. Test the model on a custom dialogue.

<i>‡</i> ====================================	
[‡] □ Step 1: Imports and Setup	
 	
mport warnings	
varnings.filterwarnings('ignore')	

import pandas as pd
import torch
import numpy as np
import matplotlib.pyplot as plt
from datasets import load_dataset
from transformers import (
 pipeline, AutoTokenizer, AutoModelForSeq2SeqLM,

```
DataCollatorForSeq2Seq, TrainingArguments, Trainer
)
# Detect device
device = 0 if torch.cuda.is available() else -1
# □ Step 2: Load CNN/DailyMail
cnn dataset = load dataset("cnn dailymail", '3.0.0', split="train[:10]")
print("□ Sample Article:\n", cnn dataset[0]['article'])
print("\n□ Reference Summary:\n", cnn dataset[0]['highlights'])
# 

Step 3: Summarization Pipeline Comparison
summary outputs = {}
# T5-small (fine-tuned)
pipe t5 = pipeline('summarization', model='ubikpt/t5-small-finetuned-
cnn', device=device)
summary outputs['t5-small'] =
pipe t5(cnn dataset[0]['article'])[0]['summary text']
# BART-large-cnn
pipe bart = pipeline('summarization', model='facebook/bart-large-cnn',
device=device)
summary outputs['bart-large'] =
pipe bart(cnn dataset[0]['article'])[0]['summary text']
# Print both summaries
for model name, summary text in summary outputs.items():
  print(f"\n□ {model name.upper()} Summary:\n{summary text}")
# □ Step 4: Load SAMSum Dataset
samsum = load dataset("samsum", trust remote code=True)
dialogue len = [len(x['dialogue'].split())) for x in samsum['train']]
summary len = [len(x['summary'].split()) for x in samsum['train']]
```

```
# Plotting lengths
length df = pd.DataFrame({'Dialogue Length': dialogue len, 'Summary
Length': summary len})
length df.hist(figsize=(10, 3))
# □ Step 5: Preprocessing
model \ ckpt = 't5-small'
tokenizer = AutoTokenizer.from pretrained(model ckpt)
model =
AutoModelForSeg2SegLM.from pretrained(model ckpt).to(torch.device
("cuda" if device == 0 else "cpu"))
def tokenize(batch):
  return tokenizer(batch['dialogue'], text target=batch['summary'],
           max length=200, truncation=True, padding="max length")
samsum tokenized = samsum.map(tokenize, batched=True)
# 'X' Step 6: Training
data collator = DataCollatorForSeq2Seq(tokenizer=tokenizer,
model=model)
training args = TrainingArguments(
  output dir="train dir",
  num train epochs=2,
  per device train batch size=4,
  per device eval batch size=4,
  evaluation strategy="epoch",
  save strategy="epoch",
  weight decay=0.01,
  learning rate=2e-5,
  gradient accumulation steps=500,
  logging dir="./logs",
  logging steps=10
)
```

```
trainer = Trainer(
  model=model,
  args=training args,
  tokenizer=tokenizer,
  data collator=data collator,
  train dataset=samsum tokenized['train'],
  eval dataset=samsum tokenized['validation']
)
trainer.train()
trainer.save model("t5 samsum summarization")
# □ Step 7: Inference on Custom Dialogue
pipe = pipeline('summarization', model="t5 samsum summarization",
tokenizer=tokenizer, device=device)
custom_dialogue = """
Laxmi Kant: what work you planning to give Tom?
Juli: I was hoping to send him on a business trip first.
Laxmi Kant: Cool. Is there any suitable work for him?
Juli: He did excellent in the last quarter. I will assign a new project once
he is back.
** ** **
print("\n \ Custom Dialogue Summary:\n",
pipe(custom dialogue)[0]['summary text'])
```

Image Classification – ViT

2. Steps Summary

Step-by-Step Summary:

- 1. Load dataset of Indian food images.
- 2. Map class labels to integers (label2id, id2label).
- 3. Load image processor compatible with Vision Transformer (ViT).
- 4. Define preprocessing pipeline: Resize, crop, normalize, and convert to tensors.
- 5. Apply transformations to training dataset.
- 6. Define evaluation metric (accuracy).
- 7. Load pre-trained ViT model and configure it for classification.
- 8. Set training arguments like learning rate, batch size, etc.
- 9. Train the model using Hugging Face Trainer.
- 10. Save the trained model and create an inference pipeline.
- 11. Download an image from the internet and use the model to predict the food class.

import warnings warnings.filterwarnings('ignore')

import requests from io import BytesIO from PIL import Image import torch import numpy as np import pandas as pd import evaluate

from datasets import load_dataset from torchvision.transforms import Compose, RandomResizedCrop, ToTensor, Normalize

from transformers import (

```
AutoImageProcessor, AutoModelForImageClassification,
  TrainingArguments, Trainer, pipeline
)
#  Step 2: Load Dataset and Prepare Labels
dataset = load dataset("rajistics/indian food images")
labels = dataset['train'].features['label'].names
label2id = {label: i for i, label in enumerate(labels)}
id2label = {i: label for i, label in enumerate(labels)}
# 

Step 3: Image Processor & Transformations
model ckpt = "google/vit-base-patch16-224-in21k"
image processor = AutoImageProcessor.from pretrained(model ckpt)
normalize = Normalize(mean=image processor.image mean,
std=image processor.image std)
size = image processor.size.get('shortest edge') or
(image processor.size['height'], image processor.size['width'])
transforms = Compose([
  RandomResizedCrop(size),
  ToTensor(),
  normalize
])
def transform fn(batch):
  batch['pixel values'] = [ transforms(img.convert("RGB")) for img in
batch['image']]
  del batch['image']
  return batch
dataset = dataset.with transform(transform fn)
# □ Step 4: Define Evaluation Metric
```

```
accuracy = evaluate.load('accuracy')
def compute metrics(eval pred):
  logits, labels = eval pred
  preds = np.argmax(logits, axis=1)
  return accuracy.compute(predictions=preds, references=labels)
# 

Step 5: Load Pretrained Model
device = torch.device("cuda" if torch.cuda.is available() else "cpu")
model = AutoModelForImageClassification.from pretrained(
  model ckpt,
  num labels=len(labels),
  id2label=id2label,
  label2id=label2id
).to(device)
# Step 6: Set Training Arguments & Train
training args = TrainingArguments(
  output dir="train dir",
  remove_unused_columns=False,
  evaluation strategy="epoch",
  save strategy="epoch",
  learning rate=5e-5,
  per device train batch size=16,
  per device eval batch size=16,
  gradient accumulation steps=4,
  num train epochs=4,
  load best model at end=True,
  metric for best model='accuracy'
)
trainer = Trainer(
  model=model,
  args=training args,
  train dataset=dataset['train'],
```

Coding_Fine_Tuning_LLM_Phi2_on_Custom_Data

```
!pip install -q accelerate -U
!pip install -q bitsandbytes -U
!pip install -q trl -U
!pip install -q peft -U
!pip install -q transformers -U
!pip install -q datasets -U
```

These install or update required libraries:

- accelerate: Efficiently trains large models with multi-GPU or mixed precision.
- bitsandbytes: Enables 8-bit optimizers and model quantization.
- trl: Hugging Face library for training RLHF and other fine-tuning.
- peft: For Parameter-Efficient Fine-Tuning (like LoRA).
- transformers: Core Hugging Face library for models, tokenizers, pipelines.
- datasets: For loading and processing datasets.

Data Handling

```
python

import pandas as pd

from datasets import load_dataset, Dataset, DatasetDict
```

- pandas: Read and manipulate the CSV file.
- load_dataset: Loads Hugging Face or custom datasets.
- Dataset and DatasetDict: Convert pandas DataFrame into HF Dataset and split into train/test.

✓ Preprocessing

Extracts the most specific category (last subcategory) from a pipe-separated string.

- · Concatenates product names and descriptions.
- Converts to HF dataset.
- · Shuffles and splits into train/test.

Formatting Function

```
python

def formatting_func(example):
    return f"""..."""
```

· Formats prompts in an instruction-following way.

Model and Tokenizer (Phi-2)

- · Loads the base Phi-2 model and tokenizer.
- · Configured for causal language modeling.

Tokenization

```
python

def tokenize(prompt): ...
```

- Tokenizes the prompt.
- · Sets input and labels (for causal LM training).

Evaluation

Generates output for a test prompt.

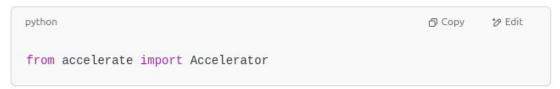
LoRA / PEFT Setup

- · Creates LoRA adapters for efficient fine-tuning.
- · Injects trainable low-rank adapters into the model.

Model Parameter Summary

· Calculates number and percentage of trainable parameters.

Accelerator



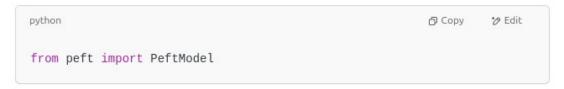
Wraps model for efficient training (handles device placement, mixed precision).

Training Setup



- TrainingArguments: Defines how training should run.
- Trainer: Manages training loop.
- DataCollatorForLanguageModeling: Collates batches and pads inputs.

Final Evaluation



Loads fine-tuned LoRA adapter weights into the base model for inference.

Zip Save

Compresses and saves the trained adapter checkpoint.

2. Summary of Fine-Tuning Steps

- 1. Install dependencies for model training, quantization, and dataset handling.
- 2. Load product data from CSV, clean and reformat.
- Format for instruction tuning standardize prompts for both tasks.
- 4. Load model/tokenizer load Phi-2 with 8-bit quantization for memory efficiency.
- 5. Tokenize and prepare dataset left padding, max length, labels for causal LM.
- 6. Evaluate base model generate from raw Phi-2 before fine-tuning.
- 7. Add LoRA adapters (QLoRA) only train a small number of parameters efficiently.
- 8. Wrap with accelerator prepare model for training on GPU/mixed precision.
- 9. Configure trainer training loop, logging, evaluation, checkpoints.
- 10. Train model fine-tune with custom formatted prompts.
- 11. Evaluate fine-tuned model load trained model, generate outputs.
- 12. Export model adapter zip adapter weights for sharing/deployment.

```
# === Install Requirements ====
!pip install -q accelerate -U
!pip install -q bitsandbytes -U
!pip install -q trl -U
!pip install -q peft -U
!pip install -q transformers -U
!pip install -q datasets -U
# === Imports ===
import torch
import pandas as pd
from datasets import Dataset
from transformers import (
  AutoTokenizer, AutoModelForCausalLM, TrainingArguments,
  Trainer, DataCollatorForLanguageModeling
from peft import LoraConfig, get peft model, PeftModel
from accelerate import Accelerator
# === Load and Prepare Dataset ===
df = pd.read csv(
  'https://github.com/laxmimerit/All-CSV-ML-Data-Files-
Download/raw/master/amazon_product_details.csv',
  usecols=['category', 'about product', 'product name']
df['category'] = df['category'].apply(lambda x: x.split('|')[-1])
products = df[['category', 'product name']].rename(columns={'product name': 'text'})
```

```
description = df[['category', 'about product']].rename(columns={'about product':
'text'})
products['task type'] = 'Product Name'
description['task type'] = 'Product Description'
df = pd.concat([products, description], ignore index=True)
dataset = Dataset.from pandas(df).shuffle(seed=0).train test split(test size=0.1)
# === Prompt Formatting ===
def formatting func(example):
  return f"""Given the product category, you need to generate a
'{example['task type']}'.
### Category: {example['category']}
### {example['task type']}: {example['text']}"""
# === Tokenization ===
base model id = "microsoft/phi-2"
tokenizer = AutoTokenizer.from pretrained(base model id, use fast=False,
padding side='left')
tokenizer.pad token = tokenizer.eos token
max length = 400
def tokenize(example):
  enc = tokenizer(formatting func(example), truncation=True,
max length=max length, padding="max length")
  enc['labels'] = enc['input ids'].copy()
  return enc
dataset = dataset.map(tokenize)
# === Load Model with LoRA ===
model = AutoModelForCausalLM.from pretrained(base model id,
torch dtype=torch.float16, load in 8bit=True)
lora config = LoraConfig(
  r=32, lora alpha=64, target modules=["Wqkv", "fc1", "fc2"],
  bias="none", lora dropout=0.05, task type="CAUSAL LM"
model = get peft model(model, lora config)
# === Print Trainable Parameters ===
def print trainable parameters(model):
  trainable = sum(p.numel() for p in model.parameters() if p.requires grad)
  total = sum(p.numel() for p in model.parameters())
  print(f"Trainable: {trainable} / {total} ({100 * trainable / total:.2f}%)")
print trainable parameters(model)
# === Accelerate & Prepare ===
accelerator = Accelerator()
model = accelerator.prepare model(model)
# === Training Setup ===
```

```
args = TrainingArguments(
  output dir="train-dir", per device train batch size=2,
  gradient accumulation steps=1, max steps=500, learning rate=2.5e-5,
  optim="paged adamw 8bit", logging steps=25, save steps=25,
  evaluation strategy="steps", eval steps=25, do eval=True
)
trainer = Trainer(
  model=model,
  args=args,
  train dataset=dataset['train'],
  eval dataset=dataset['test'],
  data collator=DataCollatorForLanguageModeling(tokenizer, mlm=False),
model.config.use cache = False
trainer.train()
# === Evaluation After Fine-Tuning ===
base model = AutoModelForCausalLM.from pretrained(base model id,
load in 8bit=True, torch dtype=torch.float16)
eval tokenizer = AutoTokenizer.from pretrained(base model id, use fast=False)
eval tokenizer.pad token = eval tokenizer.eos token
ft model = PeftModel.from pretrained(base model, './train-dir/checkpoint-500')
eval prompt = """Given the product category, you need to generate a 'Product
Description'.
### Category: BatteryChargers
### Product Description:"""
model input = eval tokenizer(eval prompt, return tensors='pt').to('cuda')
ft model.eval()
with torch.no grad():
  output = ft model.generate(**model input, max new tokens=256,
repetition penalty=1.15)
  print(eval tokenizer.decode(output[0], skip special tokens=True))
# === Save the Fine-Tuned Adapter ===
!zip -r phi2 glora adapter.zip ./train-dir/checkpoint-500
```

Fine tuning LLAMA Instruction tuning

1. Explanation of All Libraries and Functions Used

datasets

load_dataset: Loads datasets from the Hugging Face Hub. Here, it loads
 "HuggingFaceH4/ultrachat_200k".

transformers

- AutoTokenizer: Loads a tokenizer that converts between text and token IDs.
- AutoModelForCausalLM: Loads a causal language model (used for text generation).
- BitsAndBytesConfig: Enables 4-bit quantization for efficient training with QLoRA.
- TrainingArguments: Specifies training hyperparameters.
- Trainer & SFTTrainer: Trainer classes that run training loops. SFTTrainer is a special class from trl for Supervised Fine-Tuning (SFT).
- pipeline: Simplifies inference using models for tasks like text-generation.

peft (Parameter-Efficient Fine-Tuning)

- LoraConfig: Configuration for Low-Rank Adaptation (LoRA).
- get_peft_model: Wraps the base model with trainable LoRA adapters.
- prepare_model_for_kbit_training: Prepares the model for 4-bit training.
- AutoPeftModelForCausalLM: Loads a PEFT model and merges adapters for inference.

torch

Used for tensor operations and accessing GPU.

2. Summary of Fine-Tuning Steps

- 1. Load Dataset: Use UltraChat and keep 10K samples.
- Format Prompts: Convert conversations into prompt-response pairs using <|user|> and <|assistant|>.
- 3. Load Base Model & Tokenizer: Use TinyLlama-1.1B, with 4-bit quantization (nf4).
- 4. Prepare Model for QLoRA:
 - · Apply LoRA adapters with LoraConfig.
 - Target modules: q_proj, k_proj, v_proj, etc.
- 5. Set Training Arguments: Configure hyperparameters (batch size, learning rate, etc.).
- 6. Fine-tune with SFTTrainer:
 - Use the formatted text dataset.
 - · Train model with LoRA adapters only.
- 7. Merge Adapters for Inference: Merge adapters into the base model for deployment.
- 8. Generate Text: Use pipeline for evaluation on a prompt.
- 9. Export: Save adapter and zip the fine-tuned model directory.

```
# ====== INSTALL DEPENDENCIES
_____
#!pip install -q accelerate bitsandbytes trl peft transformers datasets
import torch
from datasets import load dataset
from transformers import (
  AutoTokenizer, AutoModelForCausalLM, BitsAndBytesConfig,
  TrainingArguments, pipeline
from peft import (
  LoraConfig, prepare model for kbit training,
  get peft model, AutoPeftModelForCausalLM
from trl import SFTTrainer
          ===== LOAD DATASET ===
dataset = load dataset("HuggingFaceH4/ultrachat 200k",
trust_remote_code=True, split="train_sft")
dataset = dataset.shuffle(seed=0).select(range(10 000))
```

```
# ====== PROMPT FORMATTING
tokenizer = AutoTokenizer.from pretrained("TinyLlama/TinyLlama-
1.1B-Chat-v1.0")
def format prompt(example):
  chat = example['messages']
  prompt = tokenizer.apply chat template(chat, tokenize=False)
  return {'text': prompt}
dataset = dataset.map(format_prompt)
# ===== LOAD & PREPARE BASE MODEL
model name = "TinyLlama/TinyLlama-1.1B-intermediate-step-1431k-
3T"
bnb config = BitsAndBytesConfig(
  load in 4bit=True,
  bnb 4bit quant type="nf4",
  bnb 4bit compute dtype=torch.float16,
  bnb 4bit use double quant=True
)
tokenizer = AutoTokenizer.from pretrained(model name,
trust remote code=True)
tokenizer.pad token = "<PAD>"
tokenizer.padding side = "left"
model = AutoModelForCausalLM.from pretrained(
  model name,
  device map="auto",
  quantization config=bnb config
model.config.use cache = False
model.config.pretraining tp = 1
# ====== APPLY LORA ADAPTERS
peft config = LoraConfig(
  lora alpha=32,
  lora_dropout=0.1,
  r=64,
```

```
bias="none",
  task type="CAUSAL LM",
  target_modules=["q_proj", "k_proj", "v_proj", "o_proj", "gate_proj",
"up proj", "down proj"]
model = prepare model for kbit training(model)
model = get peft model(model, peft config)
# ====== TRAINING CONFIG
training args = TrainingArguments(
  output dir="train dir",
  per device train batch size=2,
  gradient accumulation steps=4,
  optim="paged adamw 32bit",
  learning rate=2e-4,
  Ir scheduler type="cosine",
  num train epochs=1,
  logging_steps=10,
  fp16=True,
  gradient checkpointing=True
# ====== FINE-TUNE MODEL
trainer = SFTTrainer(
  model=model,
  train dataset=dataset,
  dataset text field="text",
  tokenizer=tokenizer.
  args=training args,
  max seq length=512,
  peft config=peft config
)
trainer.train()
# ====== SAVE THE ADAPTER
trainer.model.save pretrained("TinyLlama-1.1B-qlora")
```

====== LOAD MERGED MODEL FOR
INFERENCE ===================================
model = AutoPeftModelForCausalLM.from_pretrained("TinyLlama-1.1B-qlora", device map="auto")
merged_model = model.merge_and_unload()
====== INFERENCE ========
<pre>pipe = pipeline("text-generation", model=merged_model,</pre>
tokenizer=tokenizer)
prompt = """< user >\nTell me something about Large Language
Models.\n< assistant >\n"""
output = pipe(prompt, max_new_tokens=100) print(output[0]["generated text"])
print(output[o][generated_text])
====== ZIP THE CHECKPOINT
!zip -r tiny_llama_qlora_adapter.zip TinyLlama-1.1B-qlora