

✓ 1. Load Dataset:

```
1 !nvidia-smi
```

↗ Tue May 13 06:00:11 2025

NVIDIA-SMI 550.54.15				Driver Version: 550.54.15				CUDA Ver	
GPU	Name		Persistence-M	Bus-Id	Disp.A	Volati			
Fan	Temp	Perf	Pwr:Usage/Cap		Memory-Usage	GPU-Ut			
=====							=====		
0	Tesla T4		Off	00000000:00:04.0	Off				
N/A	70C	P0	30W / 70W	2466MiB / 15360MiB		0			

Processes:						
GPU	GI	CI	PID	Type	Process name	
	ID	ID				
=====						

```
1 !pip install transformers torch scikit-learn pandas
```

```

1 import re
2 import os
3 import json
4 import pandas as pd
5 import matplotlib.pyplot as plt
6 import torch
7 from sklearn.model_selection import train_test_split
8 from sklearn.metrics import accuracy_score, precision_score, recall_score,
9 #Pre-train BERT:
10 from transformers import BertTokenizer, BertForSequenceClassification, Trai
11 from sklearn.preprocessing import LabelEncoder
12 from peft import get_peft_model, LoraConfig
13 #Confidence Score System:
14 import requests
15 from Bio import Entrez
16 from langchain import LLMChain, PromptTemplate
17 from langchain.chains import RetrievalQA
18 from langchain.vectorstores import FAISS
19 from langchain.embeddings import OpenAIEmbeddings
20 #Gemini Model:
21 from google import genai
22 from google.genai import types
23 import base64
24 import google.generativeai as genai

```

(1) Covid Fake News Dataset:

```

1 # List of JSON files to process
2 json_files = [
3     'Cleaned_Covid19_Train.json',
4     'Cleaned_Covid19_Dev.json',
5 ]
6 data_dict = {}
7 # Process each JSON file
8 for json_file in json_files:
9     # Load the dataset
10    with open(json_file, 'r') as file:
11        data = json.load(file)
12
13    # Prepare a list to hold the processed data
14    jsonl_data = []
15
16    # Extract and process each entry
17    for entry in data:
18        # Extract the id, tweet, and label
19        tweet = entry['tweet']
20        label = entry['label']
21
22        # Tokenize the tweet

```

```

23     tokens = re.findall(r'\b\w+\b', tweet) # Keep only words and numbers
24     reconstructed_tweet = ' '.join(tokens)
25
26     # Prepare the JSONL entry with the required structure
27     jsonl_entry = {
28         "systemInstruction": {
29             "role": "assistant", # Example role, adjust as needed
30             "parts": [
31                 {
32                     "text": "Classification the content is Fake, Real,
33                 }
34             ]
35         },
36         "contents": [
37             {
38                 "role": "user",
39                 "parts": [
40                     {
41                         "text": f"TRANSCRIPT: \n{reconstructed_tweet}\n"
42                     }
43                 ]
44             },
45             {
46                 "role": "model",
47                 "parts": [
48                     {
49                         "text": label # The label indicating the model's prediction
50                     }
51                 ]
52             }
53         ]
54     }
55     jsonl_data.append(jsonl_entry)
56
57 # Write the processed data to a JSONL file
58 output_file = json_file.replace('.json', '.jsonl') # Change the extension
59 with open(output_file, 'w') as outfile:
60     for entry in jsonl_data:
61         json.dump(entry, outfile)
62         outfile.write('\n') # Write each entry on a new line
63     print(f"Processed {json_file} and saved to {output_file}.")
64     data_dict[json_file] = jsonl_data
65 # Access the data using the correct keys - the original filenames
66 covid_train_data = data_dict['Cleaned_Covid19_Train.json'] # Corrected key
67 covid_dev_data = data_dict['Cleaned_Covid19_Dev.json'] # Corrected key
68 # Print the first few entries for verification
69 print(f"First few entries from claims_test_data:\n{covid_train_data[:5]}")
70

```

➡ Processed Cleaned_Covid19_Train.json and saved to Cleaned_Covid19_Train.json
Processed Cleaned_Covid19_Dev.json and saved to Cleaned_Covid19_Dev.jsonl.
First few entries from claims_test_data:
[{'systemInstruction': {'role': 'assistant', 'parts': [{'text': 'Classifica

(2) Health Fact Dataset:

```
1 import json
2 import re
3 import os
4
5 # List of JSON files to process
6 json_files = [
7     'healthfact_traindata.json',
8     'cleaned_healthfact_test.json',
9     'cleaned_healthfact_dev.json'
10 ]
11 data_dict = {}
12 # Process each JSON file
13 for json_file in json_files:
14     # Prepare a list to hold the processed data
15     jsonl_data = []
16     # Load the dataset
17     with open(json_file, 'r') as file:
18         # Read each line as a separate JSON object
19         for line in file:
20             try:
21                 entry = json.loads(line)
22                 # Extract the claim, explanation, and label
23                 claim = entry['claim']
24                 explanation = entry['explanation']
25                 label = entry['label']
26
27                 # Tokenize the claim
28                 tokens = re.findall(r'\b\w+\b', claim) # Keep only words and
29                 reconstructed_claim = ' '.join(tokens)
30
31                 # Prepare the JSONL entry in the required format
32                 jsonl_entry = {
33                     "systemInstruction": {
34                         "role": "assistant", # Example role, adjust as needed
35                         "parts": [
36                             {
37                                 "text": "You are a helpful assistant." # Example text
38                             }
39                         ]
40                     },
41                     "contents": [
42                         {
```

```

43         "role": "user",
44         "parts": [
45             {
46                 "text": f"CLAIM: {reconstructed_claim}\n"
47             }
48         ]
49     },
50     {
51         "role": "model",
52         "parts": [
53             {
54                 "text": label # The label indicating t
55             }
56         ]
57     }
58 ]
59 }
60 jsonl_data.append(jsonl_entry)
61 except json.JSONDecodeError as e:
62     print(f"Error decoding JSON: {e}")
63 # Use the correct key to store the data in the dictionary – keep the or
64 data_dict[json_file] = jsonl_data
65 # Access the data using the correct keys – the original filenames
66 healthfact_train_data = data_dict['healthfact_traindata.json'] # Correctec
67 healthfact_test_data = data_dict['cleaned_healthfact_test.json'] # Correct
68 healthfact_dev_data = data_dict['cleaned_healthfact_dev.json'] # Correctec
69 # Print the first few entries for verification
70 print(f"First few entries from healthfact_train_data:\n{healthfact_train_da
71 # Optionally, write the processed data to JSONL files
72 for json_file, jsonl_data in data_dict.items():
73     output_file = json_file.replace('.json', '.jsonl') # Change the extens
74     with open(output_file, 'w') as outfile:
75         for entry in jsonl_data:
76             json.dump(entry, outfile)
77             outfile.write('\n') # Write each entry on a new line
78     print(f"Processed {json_file} and saved to {output_file}.")

```

➡ First few entries from healthfact_train_data:

```

[{'systemInstruction': {'role': 'assistant', 'parts': [{'text': 'You are a
Processed healthfact_traindata.json and saved to healthfact_traindata.jsonl
Processed cleaned_healthfact_test.json and saved to cleaned_healthfact_test
Processed cleaned_healthfact_dev.json and saved to cleaned_healthfact_dev.j

```

(3) Scifact Dataset:

```

1 import json
2 import re
3
4 # List of JSONL files to process

```

```

5 jsonl_files = [
6     'dev_3class.jsonl',
7     'train_3class.jsonl'
8 ]
9 data_dict = {}
10 # Process each JSONL file
11 for jsonl_file in jsonl_files:
12     # Prepare a list to hold the processed data
13     processed_data = []
14
15     # Load the dataset
16     with open(jsonl_file, 'r') as file:
17         for line in file:
18             try:
19                 entry = json.loads(line)
20
21                 # Extract the claim, explanation, and label
22                 claim = entry['claim']
23                 explanation = entry['evidence_text']
24                 label = entry['label']
25
26                 # Tokenize the claim
27                 tokens = re.findall(r'\b\w+\b', claim) # Keep only words and
28                 reconstructed_claim = ' '.join(tokens)
29
30                 # Prepare the JSONL entry in the required format
31                 jsonl_entry = {
32                     "systemInstruction": {
33                         "role": "assistant", # Example role, adjust as needed
34                         "parts": [
35                             {
36                                 "text": "You are a helpful assistant." # Example text
37                             }
38                         ]
39                     },
40                     "contents": [
41                         {
42                             "role": "user",
43                             "parts": [
44                                 {
45                                     "text": f"CLAIM: {reconstructed_claim}\n"
46                                 }
47                             ]
48                         },
49                         {
50                             "role": "model",
51                             "parts": [
52                                 {
53                                     "text": label # The label indicating the model's response
54                                 }
55                             ]
56                         }
57                     ]
58                 }
59                 processed_data.append(jsonl_entry)
60
61     data_dict[jsonl_file] = processed_data
62
63 # Save the data dictionary to a JSON file
64 with open('data_dict.json', 'w') as f:
65     json.dump(data_dict, f)
66
67 # Print the number of processed entries
68 print(f"Total processed entries: {len(data_dict['train_3class.jsonl']) + len(data_dict['dev_3class.jsonl'])}")

```

```

56         }
57     ]
58 }
59     # Append the modified entry to the processed data list
60     processed_data.append(jsonl_entry) # Append the processed
61 except json.JSONDecodeError as e:
62     print(f"Error decoding JSON: {e}")
63 # Store the processed data in the dictionary
64 data_dict[jsonl_file] = processed_data
65 # Access the data using the correct keys - the original filenames
66 scifact_train_data = data_dict['train_3class.jsonl'] # Corrected key
67 scifact_test_data = data_dict['dev_3class.jsonl'] # Corrected key
68 # Print the first few entries for verification
69 print(f"First few entries from scifact_train_data:\n{scifact_train_data[:5]}")
70 # Optionally, write the processed data to new JSONL files
71 for jsonl_file, processed_data in data_dict.items():
72     output_file = jsonl_file.replace('.jsonl', '_processed.jsonl') # Change
73     with open(output_file, 'w') as outfile:
74         for entry in processed_data:
75             json.dump(entry, outfile)
76             outfile.write('\n') # Write each entry on a new line
77     print(f"Processed {jsonl_file} and saved to {output_file}.")

```

➡ First few entries from scifact_train_data:
 [{"systemInstruction": {"role": "assistant", "parts": [{"text": "You are a
 Processed dev_3class.jsonl and saved to dev_3class_processed.jsonl.
 Processed train_3class.jsonl and saved to train_3class_processed.jsonl.

✓ 2. Data Exploration

```

1 # Function to explore a dataset
2 def explore_dataset(data, dataset_name):
3     print(f"Exploring dataset: {dataset_name}")
4     print(f"Number of entries: {len(data)}")
5
6     # Convert to DataFrame for easier analysis
7     df = pd.DataFrame(data)
8
9     # Display the first few entries
10    print("First few entries:")
11    print(df.head())
12
13    # Display basic statistics
14    print("\nBasic statistics:")
15    print(df.describe(include='all'))
16
17    # Check the distribution of labels (if applicable)
18    if 'label' in df.columns:

```

```

19     label_counts = df['label'].value_counts()
20     print("\nLabel distribution:")
21     print(label_counts)
22
23     # Plot the label distribution
24     label_counts.plot(kind='bar', title='Label Distribution')
25     plt.xlabel('Labels')
26     plt.ylabel('Counts')
27     plt.show()
28
29     print("\n" + "-" * 40 + "\n")
30
31 # Explore each dataset
32 explore_dataset(covid_train_data, "Cleaned Covid19 Train Data")
33 explore_dataset(healthfact_train_data, "Healthfact Train Data")
34 explore_dataset(scifact_train_data, "SciFact Train Data")

```

```

0 {'role': 'assistant', 'parts': [{'text': 'You ...
1 {'role': 'assistant', 'parts': [{'text': 'You ...
2 {'role': 'assistant', 'parts': [{'text': 'You ...
3 {'role': 'assistant', 'parts': [{'text': 'You ...
4 {'role': 'assistant', 'parts': [{'text': 'You ...

```

```

                                contents
0  [{'role': 'user', 'parts': [{'text': 'CLAIM: T...
1  [{'role': 'user', 'parts': [{'text': 'CLAIM: A...
2  [{'role': 'user', 'parts': [{'text': 'CLAIM: S...
3  [{'role': 'user', 'parts': [{'text': 'CLAIM: S...
4  [{'role': 'user', 'parts': [{'text': 'CLAIM: S...

```

Basic statistics:

```

                                systemInstruction \
count                                9804
unique                                1
top      {'role': 'assistant', 'parts': [{'text': 'You ...
freq                                9804

```

```

                                contents
count                                9804
unique                                9803
top      [{'role': 'user', 'parts': [{'text': 'CLAIM: P...
freq                                2

```

Exploring dataset: SciFact Train Data

Number of entries: 1261

First few entries:

```

                                systemInstruction \
0  {'role': 'assistant', 'parts': [{'text': 'You ...
1  {'role': 'assistant', 'parts': [{'text': 'You ...
2  {'role': 'assistant', 'parts': [{'text': 'You ...
3  {'role': 'assistant', 'parts': [{'text': 'You ...
4  {'role': 'assistant', 'parts': [{'text': 'You ...

```

contents


```

0  [{'role': 'user', 'parts': [{'text': 'CLAIM: 0...
1  [{'role': 'user', 'parts': [{'text': 'CLAIM: 1...
2  [{'role': 'user', 'parts': [{'text': 'CLAIM: 1...
3  [{'role': 'user', 'parts': [{'text': 'CLAIM: 1...
4  [{'role': 'user', 'parts': [{'text': 'CLAIM: 3...

```

Basic statistics:

	systemInstruction \
count	1261
unique	1
top	{'role': 'assistant', 'parts': [{'text': 'You ...
freq	1261

	contents
count	1261
unique	1258
top	[{'role': 'user', 'parts': [{'text': 'CLAIM: A...
freq	2

✓ 3. Training Strategy

```

1 # Convert datasets to DataFrames for easier manipulation
2 healthfact_df = pd.DataFrame(healthfact_train_data)
3 scifact_df = pd.DataFrame(scifact_train_data)
4
5 # Combine HealthFact and SciFact datasets for pre-training
6 combined_pretrain_df = pd.concat([healthfact_df, scifact_df], ignore_index=
7
8 # Save the combined dataset for pre-training
9 combined_pretrain_df.to_json('combined_pretrain_data.jsonl', orient='recorc
10
11 # Convert COVID-19 dataset to DataFrame
12 covid_df = pd.DataFrame(covid_train_data)
13
14 # Save the COVID-19 dataset for fine-tuning
15 covid_df.to_json('covid_finetune_data.jsonl', orient='records', lines=True)
16
17 print("Datasets combined and saved for train dataset:")
18 print("1. Combined Pre-train Data: combined_pretrain_data.jsonl")
19 print("2. COVID-19 Fine-tune Data: covid_finetune_data.jsonl")

```



Datasets combined and saved for train dataset:

1. Combined Pre-train Data: combined_pretrain_data.jsonl
2. COVID-19 Fine-tune Data: covid_finetune_data.jsonl

```

1 # Convert datasets to DataFrames for easier manipulation
2 healthfact_df_test = pd.DataFrame(healthfact_test_data)
3 scifact_df_test = pd.DataFrame(scifact_test_data)
4
5 # Combine HealthFact and SciFact datasets for pre-training
6 combined_pretrain_df_test = pd.concat([healthfact_df_test, scifact_df_test]
7
8 # Save the combined dataset for pre-training
9 combined_pretrain_df_test.to_json('combined_pretrain_test_data.jsonl', orie
10
11 # Convert COVID-19 dataset to DataFrame
12 covid_df_test = pd.DataFrame(covid_dev_data)
13
14 # Save the COVID-19 dataset for fine-tuning
15 covid_df_test.to_json('covid_finetune_test_data.jsonl', orient='records', 1
16
17 print("Datasets combined and saved for Test dataset:")
18 print("1. Combined Pre-train Data: combined_pretrain_test_data.jsonl")
19 print("2. COVID-19 Fine-tune Data: covid_finetune_test_data.jsonl")

```

⇒ Datasets combined and saved for Test dataset:

1. Combined Pre-train Data: combined_pretrain_test_data.jsonl
2. COVID-19 Fine-tune Data: covid_finetune_test_data.jsonl

✓ 4. Before Fine Tuning Gemini 2.0 Flash Model Prompt + Label

```

1 # Define a fine-tuning function using Gemini API
2 def generate_response(prompt):
3     response = model.generate_content(prompt)
4     return response.text
5
6 # Example few-shot training prompt
7 prompt = """
8 Claim: "6 10 Sky s EdConwaySky explains the latest COVID19 data and governn
9 """
10
11 response = generate_response(prompt)
12 print(response)

```

This claim appears to be a tweet or social media post promoting a segment on Sky News with Ed Conway explaining the latest COVID-19 data and government announcement. It also provides links to further information.

Here's a breakdown of the elements:

- **"6 10"**: This likely refers to the time the tweet was posted, possibly 6:10 AM or PM.
- **"Sky s"**: This is likely a shortened form of "Sky News's".
- **"EdConwaySky"**: This is probably the Twitter handle for Ed Conway, who is likely a Sky News correspondent.
- **"explains the latest COVID19 data and government announcement"**: This describes the content of the segment being promoted.
- **"Get more on the coronavirus data here"**: This is a call to action, encouraging viewers to click on the provided links.
- **"https t co jvGZISbFjH https t co PygSKXesBg"**: These are shortened URLs likely leading to Sky News's website or relevant articles. It's important to note that link shorteners like "t.co" can hide the true destination of the link.

Potential Issues and Things to Consider:

- **Data Accuracy**: While the claim itself isn't making a specific factual statement, the accuracy of the COVID-19 data presented in the Sky News segment would be dependent on the sources used by Ed Conway.
- **Bias**: It's important to be aware of potential biases. Sky News, like any news organization, has a perspective. Viewers should critically evaluate the information presented.
- **Outdated Information**: COVID-19 data and government announcements change rapidly. The information presented in the segment might be outdated by the time you see the tweet.
- **Link Safety**: Always be cautious when clicking on shortened links, especially from unfamiliar sources. While Sky News is a reputable organization, it's good practice to be vigilant. You can use a link expander to see the actual URL before clicking.

In Conclusion:

The claim itself is a straightforward promotion of a news segment. However, critical evaluation of the information presented in the segment is still necessary. Consider the source, potential biases, the date of the information, and always be cautious when clicking on links.

✓ 5. Model Initialization

```
1 !pip install --upgrade google-genai
2 !gcloud auth application-default login
3 !pip install --upgrade google-cloud-aiplatform
```

```
1
2 from google.colab import auth as google_auth
3 google_auth.authenticate_user()
4
5 import vertexai
6 from vertexai.generative_models import GenerativeModel
7 from vertexai.preview.tuning import sft
8
9 vertexai.init(project="sit319-25t1-nguyen-ae806d0", location="us-central1")#
10
11 gemini_pro = GenerativeModel("gemini-2.0-flash-lite-001")
12
13 sft_tuning_job = sft.train(
14     source_model=gemini_pro,
15     train_dataset="gs://daftt/Cleaned_Covid19_Train-7.jsonl",
16     tuned_model_display_name="covid_tuning",
17     epochs=100,
18     learning_rate_multiplier=1,
19 )
20
```

[/usr/local/lib/python3.11/dist-packages/google/auth/_default.py:76](#): UserWarning: Your application has authenticated using end user credentials from Google Cloud SDK without a quota project. You might receive a "quota exceeded" or "API not enabled" error. See the following page for troubleshooting: <https://cloud.google.com/docs/authentication/adc-troubleshooting/user-creds>. warnings.warn(_CLOUD_SDK_CREDENTIALS_WARNING)

INFO:vertexai.tuning._tuning:Creating SupervisedTuningJob [/usr/local/lib/python3.11/dist-packages/google/auth/_default.py:76](#): UserWarning: Your application has authenticated using end user credentials from Google Cloud SDK without a quota project. You might receive a "quota exceeded" or "API not enabled" error. See the following page for troubleshooting: <https://cloud.google.com/docs/authentication/adc-troubleshooting/user-creds>. warnings.warn(_CLOUD_SDK_CREDENTIALS_WARNING)

INFO:vertexai.tuning._tuning:SupervisedTuningJob created. Resource name: projects/181085238689/locations/us-central1/tuningJobs/1012173862948831232

INFO:vertexai.tuning._tuning:To use this SupervisedTuningJob in another session:

INFO:vertexai.tuning._tuning:tuning_job =
sft.SupervisedTuningJob('projects/181085238689/locations/us-central1/tuningJobs/1012173862948831232')

INFO:vertexai.tuning._tuning:View Tuning Job: <https://console.cloud.google.com/vertex-ai/generative/language/locations/us-central1/tuning/tuningJob/1012173862948831232?project=181085238689>

✓ 6. Training and Evaluation fine tuning: Gemini 2.0 Flash + BERT

A. Pre-train on HealthFact + SciFact for General fact-checking ability.

```
1 import pandas as pd
2 import torch
3 from sklearn.metrics import accuracy_score, precision_score, recall_score,
4 from transformers import BertTokenizer, BertForSequenceClassification, Trai
5 from sklearn.preprocessing import LabelEncoder
6 # Load the combined pre-training dataset (HealthFact + SciFact)
7 train_combined_data = pd.read_json('combined_pretrain_data.jsonl', lines=Tr
8 val_combined_data = pd.read_json('combined_pretrain_test_data.jsonl', lines
9
10 # Load the COVID-19 fine-tuning dataset
11 train_covid_data = pd.read_json('covid_finetune_data.jsonl', lines=True)
```

```

12 val_covid_data = pd.read_json('covid_finetune_test_data.jsonl', lines=True)
13
14 # Assuming the datasets have 'claim' and 'label' columns
15 # Extract claims and labels from nested structure for pre-training
16 train_claims = train_combined_data['contents'].apply(lambda x: x[0]['parts']
17 train_labels = train_combined_data['contents'].apply(lambda x: x[1]['parts']
18 val_claims = val_combined_data['contents'].apply(lambda x: x[0]['parts'][0]
19 val_labels = val_combined_data['contents'].apply(lambda x: x[1]['parts'][0]
20
21 # Convert string labels to integers
22 label_encoder = LabelEncoder()
23 train_labels = label_encoder.fit_transform(train_labels)
24 val_labels = label_encoder.transform(val_labels)
25
26 # Load the BERT tokenizer
27 tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
28
29 # Tokenize the input data for pre-training
30 train_encodings = tokenizer(train_claims, truncation=True, padding=True, ma
31 val_encodings = tokenizer(val_claims, truncation=True, padding=True, max_le
32
33 # Create a dataset class
34 class ClaimsDataset(torch.utils.data.Dataset):
35     def __init__(self, encodings, labels):
36         self.encodings = encodings
37         self.labels = labels
38
39     def __getitem__(self, idx):
40         item = {key: torch.tensor(val[idx]) for key, val in self.encodings.
41         item['labels'] = torch.tensor(self.labels[idx])
42         return item
43
44     def __len__(self):
45         return len(self.labels)
46
47 # Create datasets for pre-training
48 train_dataset = ClaimsDataset(train_encodings, train_labels)
49 val_dataset = ClaimsDataset(val_encodings, val_labels)
50
51 # Load the BERT model
52 model = BertForSequenceClassification.from_pretrained('bert-base-uncased',
53
54 # Define training arguments for pre-training with validation loss logging
55 training_args = TrainingArguments(
56     output_dir='./results/pretrain',
57     num_train_epochs=3,
58     per_device_train_batch_size=8,
59     per_device_eval_batch_size=8,
60     warmup_steps=500,
61     weight_decay=0.01,
62     logging_dir='./logs/pretrain',

```

```

63     logging_steps=10,
64     eval_strategy="epoch", # Updated to eval_strategy
65 )
66
67 # Create a Trainer instance for pre-training
68 trainer = Trainer(
69     model=model,
70     args=training_args,
71     train_dataset=train_dataset,
72     eval_dataset=val_dataset,
73     compute_metrics=lambda p: {
74         'accuracy': accuracy_score(p.label_ids, p.predictions.argmax(-1)),
75         'precision': precision_score(p.label_ids, p.predictions.argmax(-1)),
76         'recall': recall_score(p.label_ids, p.predictions.argmax(-1), average='weighted'),
77         'f1': f1_score(p.label_ids, p.predictions.argmax(-1), average='weighted'),
78         'roc_auc': roc_auc_score(p.label_ids, torch.softmax(torch.tensor(p.predictions), dim=-1).cpu().numpy())
79     },
80 )
81
82 # Pre-train the model
83 trainer.train()
84 # Save the model and tokenizer
85 model_save_path = "./my_trained_model" # Choose your desired save path
86 model.save_pretrained(model_save_path)
87 tokenizer.save_pretrained(model_save_path)
88
89 print(f"Model and tokenizer saved to: {model_save_path}")

```

```
→ /usr/local/lib/python3.11/dist-packages/huggingface_hub/utils/_auth.py:94:
The secret `HF_TOKEN` does not exist in your Colab secrets.
To authenticate with the Hugging Face Hub, create a token in your 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
warnings.warn(
```

tokenizer_config.json: 100% 48.0/48.0 [00:00<00:00, 4.89kB/s]

vocab.txt: 100% 232k/232k [00:00<00:00, 2.60MB/s]

tokenizer.json: 100% 466k/466k [00:00<00:00, 5.64MB/s]

config.json: 100% 570/570 [00:00<00:00, 46.3kB/s]

model.safetensors: 100% 440M/440M [00:06<00:00, 102MB/s]

Some weights of BertForSequenceClassification were not initialized from the You should probably TRAIN this model on a down-stream task to be able to us

```
wandb: WARNING The `run_name` is currently set to the same value as `Traini
```

wandb: Logging into wandb.ai. (Learn how to deploy a W&B server locally: [ht](#)

wandb: You can find your API key in your browser here: <https://wandb.ai/aut>

wandb: Paste an API key from your profile and hit enter:

wandb: **WARNING** If you're specifying your api key in code, ensure this code

```
wandb: WARNING Consider setting the WANDB API KEY environment variable, or
```

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: hieunguyen23032001 (hieunguyen23032001-deaki

Tracking run with wandb version 0.19.11

Run data is saved locally in `/content/wandb/run-20250513_052637-k7ad0pp8`

Syncing run `./results/pretrain` to [Weights & Biases \(docs\)](#)

View project at <https://wandb.ai/hieunguyen23032001-deakin-university/huggingface>

View run at <https://wandb.ai/hieunguyen23032001-deakin-university/huggingface/runs/k7ad0pp8>

[4152/4152 15:29, Epoch 3/3]

Epoch	Training Loss	Validation Loss	Accuracy	Precision	Recall	F1	Roc Auc
1	0.297300	0.288650	0.875817	0.885717	0.875817	0.877511	0.989003
2	0.357900	0.281554	0.890077	0.892447	0.890077	0.889494	0.991119
3	0.208000	0.349890	0.898990	0.902847	0.898990	0.900310	0.991414

Model and tokenizer saved to: ./my trained model

B. Training and Evaluation Covid Fake News Dataset


```

1 # Prepare the training and validation data
2 # Prepare the training and validation data
3 train_covid_claims = train_covid_data['contents'].apply(lambda x: x[0]['par
4 train_covid_labels = train_covid_data['contents'].apply(lambda x: x[1]['par
5 val_covid_claims = val_covid_data['contents'].apply(lambda x: x[0]['parts']
6 val_covid_labels = val_covid_data['contents'].apply(lambda x: x[1]['parts']
7
8 # Convert string labels to integers
9 label_encoder = LabelEncoder()
10 train_covid_labels = label_encoder.fit_transform(train_covid_labels)
11 val_covid_labels = label_encoder.transform(val_covid_labels)
12
13 # Load the BERT tokenizer
14 tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
15
16 # Tokenize the input data for pre-training
17 train_encodings = tokenizer(train_covid_claims, truncation=True, padding=Tr
18 val_encodings = tokenizer(val_covid_claims, truncation=True, padding=True,
19
20 # Create datasets for pre-training
21 train_dataset = ClaimsDataset(train_encodings, train_labels)
22 val_dataset = ClaimsDataset(val_encodings, val_labels)
23
24 # Load the BERT model
25 model = BertForSequenceClassification.from_pretrained('bert-base-uncased',
26
27 # Pre-train the model
28 trainer.train()

```

➡ Some weights of BertForSequenceClassification were not initialized from the
 You should probably TRAIN this model on a down-stream task to be able to us
 [4152/4152 15:42, Epoch 3/3]

Epoch	Training Loss	Validation Loss	Accuracy	Precision	Recall	F1	Roc Auc
1	0.401200	0.396093	0.887701	0.890988	0.887701	0.888816	0.990407
2	0.177600	0.489731	0.888295	0.889070	0.888295	0.887639	0.990609
3	0.119300	0.606078	0.888889	0.892492	0.888889	0.890275	0.990721

TrainOutput(global_step=4152, training_loss=0.15090772818669906, metrics=
 {'train runtime': 942.5088, 'train samples per second': 35.22,

✓ 7. Test Prompt and Label

```

1 def generate():
2     client = genai.Client(
3         vertexai=True,
4         project="181085238680"

```

```

4     project= 181085238689 ,
5     location="us-central1",
6 )
7
8 msg3_text1 = types.Part.from_text(text="""Clearly the Obama administration
9
10 model = "projects/181085238689/locations/us-central1/endpoints/54197709897
11 contents = [
12     types.Content(
13         role="user",
14         parts=[
15             types.Part.from_text(text="""Multiple Facebook posts claim that Auss
16         ]
17     ),
18     types.Content(
19         role="model",
20         parts=[
21             types.Part.from_text(text=label)
22         ]
23     ),
24     types.Content(
25         role="user",
26         parts=[
27             msg3_text1
28         ]
29     ),
30 ]
31 generate_content_config = types.GenerateContentConfig(
32     temperature = 0.2,
33     top_p = 0.8,
34     max_output_tokens = 1024,
35     response_modalities = ["TEXT"],
36     safety_settings = [types.SafetySetting(
37         category="HARM_CATEGORY_HATE_SPEECH",
38         threshold="OFF"
39     ),types.SafetySetting(
40         category="HARM_CATEGORY_DANGEROUS_CONTENT",
41         threshold="OFF"
42     ),types.SafetySetting(
43         category="HARM_CATEGORY_SEXUALLY_EXPLICIT",
44         threshold="OFF"
45     ),types.SafetySetting(
46         category="HARM_CATEGORY_HARASSMENT",
47         threshold="OFF"
48     )],
49 )
50
51 for chunk in client.models.generate_content_stream(
52     model = model,
53     contents = contents,
54     config = generate_content_config,
55 ):

```

```
56     print(chunk.text, end='')
57
58 generate()
```

/usr/local/lib/python3.11/dist-packages/google/auth/_default.py:76: UserWarning: Your application has authenticated using end user credentials from Google Cloud SDK without a quota project. You might receive a "quota exceeded" or "API not enabled" error. See the following page for troubleshooting: <https://cloud.google.com/docs/authentication/adc-troubleshooting/user-creds>. warnings.warn(_CLOUD_SDK_CREDENTIALS_WARNING)

fake

✓ 8. Implement the Confidence Scoring System

```
1 !pip install transformers torch requests beautifulsoup4
```

 [Show hidden output](#)

```
1 !pip install biopython
```

 [Show hidden output](#)

```
1 !pip install langchain-community
```

 [Show hidden output](#)

Retrieval article

```

1 import requests
2 from bs4 import BeautifulSoup
3
4 def retrieve_articles(query):
5     """
6     Retrieve articles from PubMed based on a query.
7
8     This function uses the PubMed API to search for relevant articles
9     based on the provided query and parses the HTML response using Beautiful
10    """
11    base_url = "https://pubmed.ncbi.nlm.nih.gov/"
12    search_url = f"{base_url}?term={query}"
13
14    response = requests.get(search_url)
15
16    if response.status_code == 200:
17        soup = BeautifulSoup(response.content, 'html.parser')
18        # Extract article titles and summaries (example, you may need to ac
19        articles = []
20        for article_tag in soup.find_all('div', class_='docsum'): # Examp
21            title = article_tag.find('a', class_='docsum-title').text.strip
22            summary = article_tag.find('div', class_='abstract').text.strip
23            articles.append({'title': title, 'summary': summary})
24
25        return articles
26    else:
27        return None

```

Implement "Trust Score" calculation

```

1 def calculate_trust_score(prediction, retrieved_articles):
2     """
3     Calculate a trust score based on the LLM's prediction and the retrieved
4     The trust score is determined by the number of articles that support or
5     """
6     support_count = 0
7     contradict_count = 0
8
9     for article in retrieved_articles:
10         if prediction.lower() in article['title'].lower() or prediction.lower() in article['content'].lower():
11             support_count += 1
12         else:
13             contradict_count += 1
14
15     total_articles = support_count + contradict_count
16     if total_articles == 0:
17         return 0.0 # No articles found
18
19     trust_score = support_count / total_articles # Simple ratio of support
20     return trust_score

```

Implement LLM Responses

```

1 def predict_with_confidence(claim):
2     """
3     Predict the label for a claim and calculate the trust score based on re
4     """
5     model.eval()
6     with torch.no_grad():
7         inputs = prepare_input(claim)
8         outputs = model(**inputs)
9         logits = outputs.logits
10        predictions = torch.argmax(logits, dim=-1).item()
11
12    # Convert predicted label to word
13    predicted_label_word = label_encoder.inverse_transform([predictions])[0]
14
15    # Retrieve articles related to the claim
16    retrieved_articles = retrieve_articles(claim)
17
18    # Calculate the trust score
19    trust_score = calculate_trust_score(claim, retrieved_articles)
20
21    # Flag low-confidence responses
22    if trust_score < 0.5: # Example threshold
23        print(f"Low confidence for claim: '{claim}'. Trust score: {trust_sc
24    else:
25        print(f"High confidence for claim: '{claim}'. Trust score: {trust_s
26
27    return predicted_label_word, trust_score

```

Execute

```

1 claim = "Study Vaccine for Breast Ovarian Cancer Has Potential"
2 predicted_label_word, trust_score = predict_with_confidence(claim)
3 print(f"Predicted label for the claim '{claim}': '{predicted_label_word}',

```

➞ Low confidence for claim: 'Study Vaccine for Breast Ovarian Cancer Has Pote
 Predicted label for the claim 'Study Vaccine for Breast Ovarian Cancer Has

