# 1. Load Dataset:

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

Show hidden output

```python
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 from transformers import BertTokenizer, BertForSequenceClassification, Trai
10 from peft import get_peft_model, LoraConfig
```

(1) Covid Fake News Dataset:

```python
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 numbe
24        reconstructed_tweet = ' '.join(tokens)
25
```

```python
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
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 extens
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
64    print(f"Processed {json_file} and saved to {output_file}.")
65
66    data_dict[json_file] = jsonl_data
67 # Access the data using the correct keys - the original filenames
68 covid_train_data = data_dict['Cleaned_Covid19_Train.json']  # Corrected key
69 covid_dev_data = data_dict['Cleaned_Covid19_Dev.json']  # Corrected key
70 # Print the first few entries for verification
71 print(f"First few entries from claims_test_data:\n{covid_train_data[:5]}")
72
```

Processed Cleaned_Covid19_Train.json and saved to Cleaned_Covid19_Train.jso
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:

```python
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
12 data_dict = {}
13
14 # Process each JSON file
15 for json_file in json_files:
16     # Prepare a list to hold the processed data
17     jsonl_data = []
18
19     # Load the dataset
20     with open(json_file, 'r') as file:
21         # Read each line as a separate JSON object
22         for line in file:
23             try:
24                 entry = json.loads(line)
25                 # Extract the claim_id, claim, explanation, and label
26                 claim = entry['claim']
27                 explanation = entry['explanation']
28                 label = entry['label']
29
30                 # Tokenize the claim
31                 tokens = re.findall(r'\b\w+\b', claim)  # Keep only words a
32                 reconstructed_claim = ' '.join(tokens)
33
34                 # Prepare the JSONL entry
35                 jsonl_entry = {
36                     "claim": reconstructed_claim,
37                     "explanation": explanation,
38                     "label": label
39                 }
40                 jsonl_data.append(jsonl_entry)
41             except json.JSONDecodeError as e:
42                 print(f"Error decoding JSON: {e}")
```

```
43
44    # Use the correct key to store the data in the dictionary - keep the or
45    data_dict[json_file] = jsonl_data
46
47 # Access the data using the correct keys - the original filenames
48 healthfact_train_data = data_dict['healthfact_traindata.json']  # Corrected
49 healthfact_test_data = data_dict['cleaned_healthfact_test.json']  # Correct
50 healthfact_dev_data = data_dict['cleaned_healthfact_dev.json']  # Corrected
51 # Print the first few entries for verification
52 print(f"First few entries from claims_test_data:\n{healthfact_train_data[:5
```

First few entries from claims_test_data:
[{'claim': 'The money the Clinton Foundation took from from foreign governm


(3) Scifact Dataset:

```python
1  # List of JSONL files to process
2  jsonl_files = [
3      'dev_3class.jsonl',
4      'train_3class.jsonl'
5  ]
6  data_dict = {}
7  # Process each JSONL file
8  for jsonl_file in jsonl_files:
9      # Prepare a list to hold the processed data
10     processed_data = []
11
12     # Load the dataset
13     with open(jsonl_file, 'r') as file:
14         for line in file:
15             try:
16                 entry = json.loads(line)
17
18                 # Extract the claim_id, claim, explanation, and label
19                 claim = entry['claim']
20                 explanation = entry['evidence_text']
21                 label = entry['label']
22
23                 # Tokenize the claim
24                 tokens = re.findall(r'\b\w+\b', claim)  # Keep only words a
25                 reconstructed_claim = ' '.join(tokens)
26
27                 # Prepare the JSONL entry
28                 jsonl_entry = {
29                     "claim": reconstructed_claim,
30                     "explanation": explanation,
31                     "label": label
32                 }
33
34                 # Append the modified entry to the processed data list
35                 processed_data.append(jsonl_entry) # Changed from entry to
36             except json.JSONDecodeError as e:
37                 print(f"Error decoding JSON: {e}")
38     data_dict[jsonl_file] = processed_data
39
40  scifact_train_data = data_dict['train_3class.jsonl']  # Corrected key
41  scifact_test_data = data_dict['dev_3class.jsonl']  # Corrected key
42  print(f"First few entries from claims_test_data:\n{scifact_train_data[:5]}'
```

```
First few entries from claims_test_data:
[{'claim': '0 dimensional biomaterials lack inductive properties', 'explana
```

## ∨ 2. Data Exploration

```python
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")
```

```
Exploring dataset: Cleaned Covid19 Train Data
Number of entries: 6420
First few entries:
                                                claim label
0  The CDC currently reports 99031 deaths In gene...  real
1  States reported 1121 deaths a small rise from ...  real
2  Politically Correct Woman Almost Uses Pandemic...  fake
3  IndiaFightsCorona We have 1524 COVID testing l...  real
4  Populous states can generate large case counts...  real

Basic statistics:
                                                claim label
count                                            6420  6420
unique                                           6379     2
top     FREE HORSES 52 thoroughbred horses need homes ...  real
freq                                                3  3360

Label distribution:
```
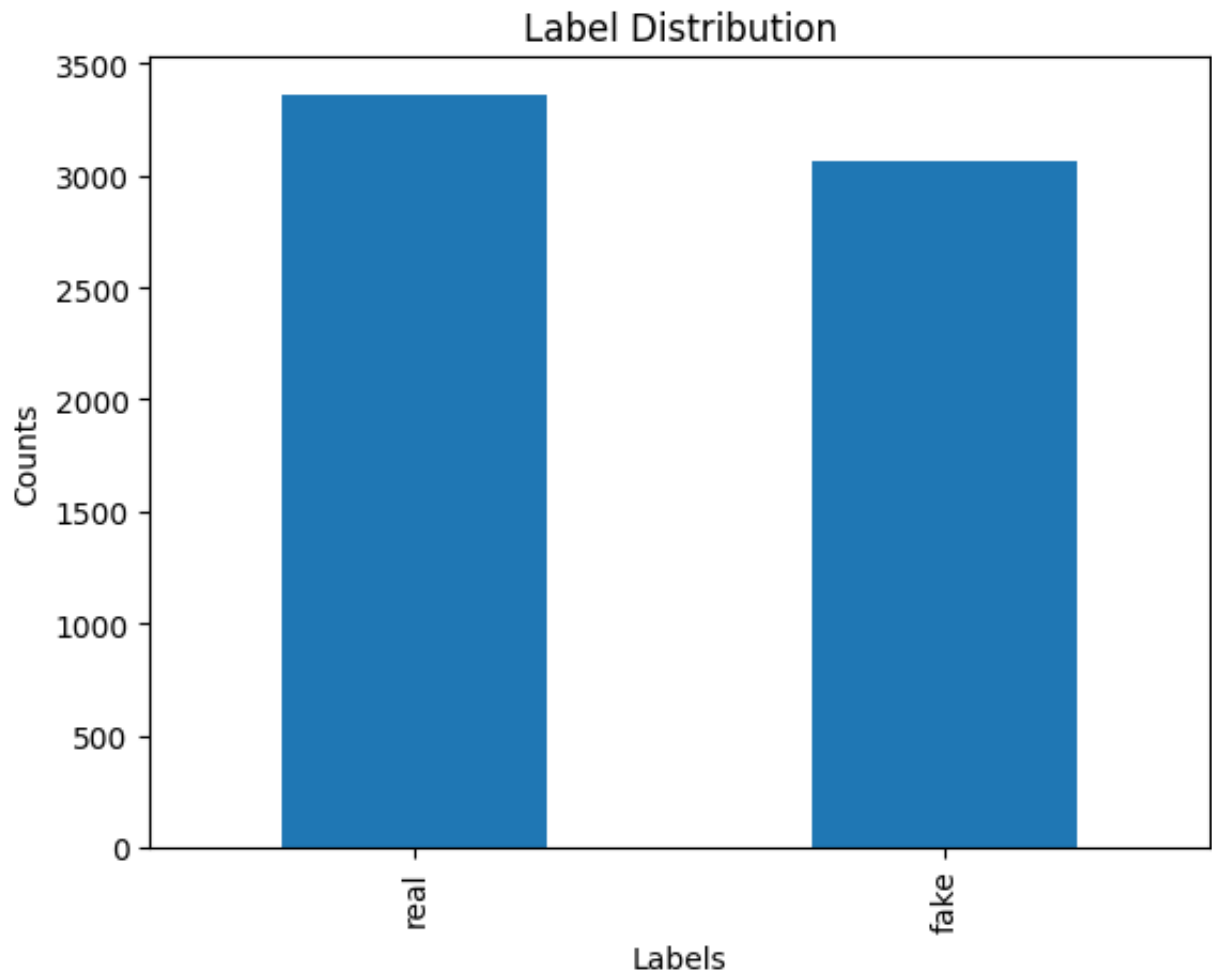
```
label
real    3360
fake    3060
Name: count, dtype: int64
```

## Label Distribution



----------------------------------------

```
Exploring dataset: Healthfact Train Data
Number of entries: 9804
First few entries:
                                            claim  \
0  The money the Clinton Foundation took from fro...
1    Annual Mammograms May Have More False Positives
2  SBRT Offers Prostate Cancer Patients High Canc...
3  Study Vaccine for Breast Ovarian Cancer Has Po...
4  Some appendicitis cases may not require emerge...

                                       explanation       label
0  "Gingrich said the Clinton Foundation ""took m...       false
1  This article reports on the results of a study...  MISLEADING
2  This news release describes five-year outcomes...  MISLEADING
3  While the story does many things well, the ove...        true
4  We really don't understand why only a handful ...        true

Basic statistics:
                                            claim  \
count                                        9804
unique                                       9799
top      Brain scans predict which dyslexics will read
freq                                            2
```

```
                                    explanation  label
count                                      9804   9804
unique                                     9660      3
top      A U.S. judge on Wednesday appointed prominent ...   true
freq                                          2   5078

Label distribution:
label
true           5078
false          3001
MISLEADING     1725
Name: count, dtype: int64
```
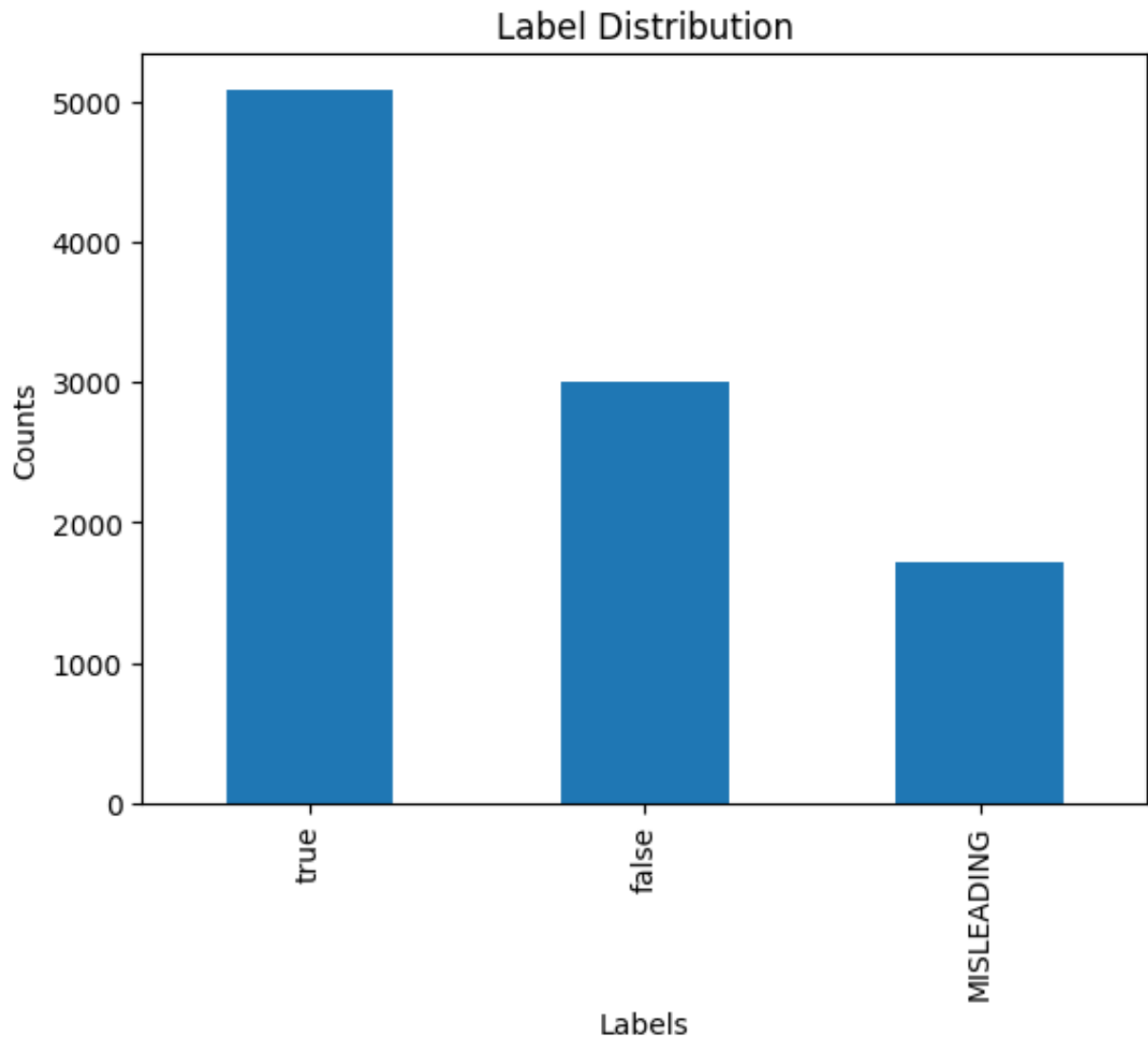
## Label Distribution



----------------------------------------

```
Exploring dataset: SciFact Train Data
Number of entries: 1261
First few entries:
                                            claim  \
0   0 dimensional biomaterials lack inductive prop...
1   1 in 5 million in UK have abnormal PrP positivity
2   1 1 of colorectal cancer patients are diagnose...
3   10 of sudden infant death syndrome SIDS deaths...
4   32 of liver transplantation programs required ...
```

```
                                    explanation      label
0                                                Misleading
1   RESULTS Of the 32,441 appendix samples 16 were...     False
2                                                Misleading
3                                                Misleading
4   Policies requiring discontinuation of methadon...      True

Basic statistics:
                                            claim explanation label
count                                        1261        1261  1261
unique                                        807         691     3
top       Energy balance requires hypothalamic glutamate...         True
freq                                           11         304   616

Label distribution:
label
True          616
False         341
Misleading    304
Name: count, dtype: int64
```
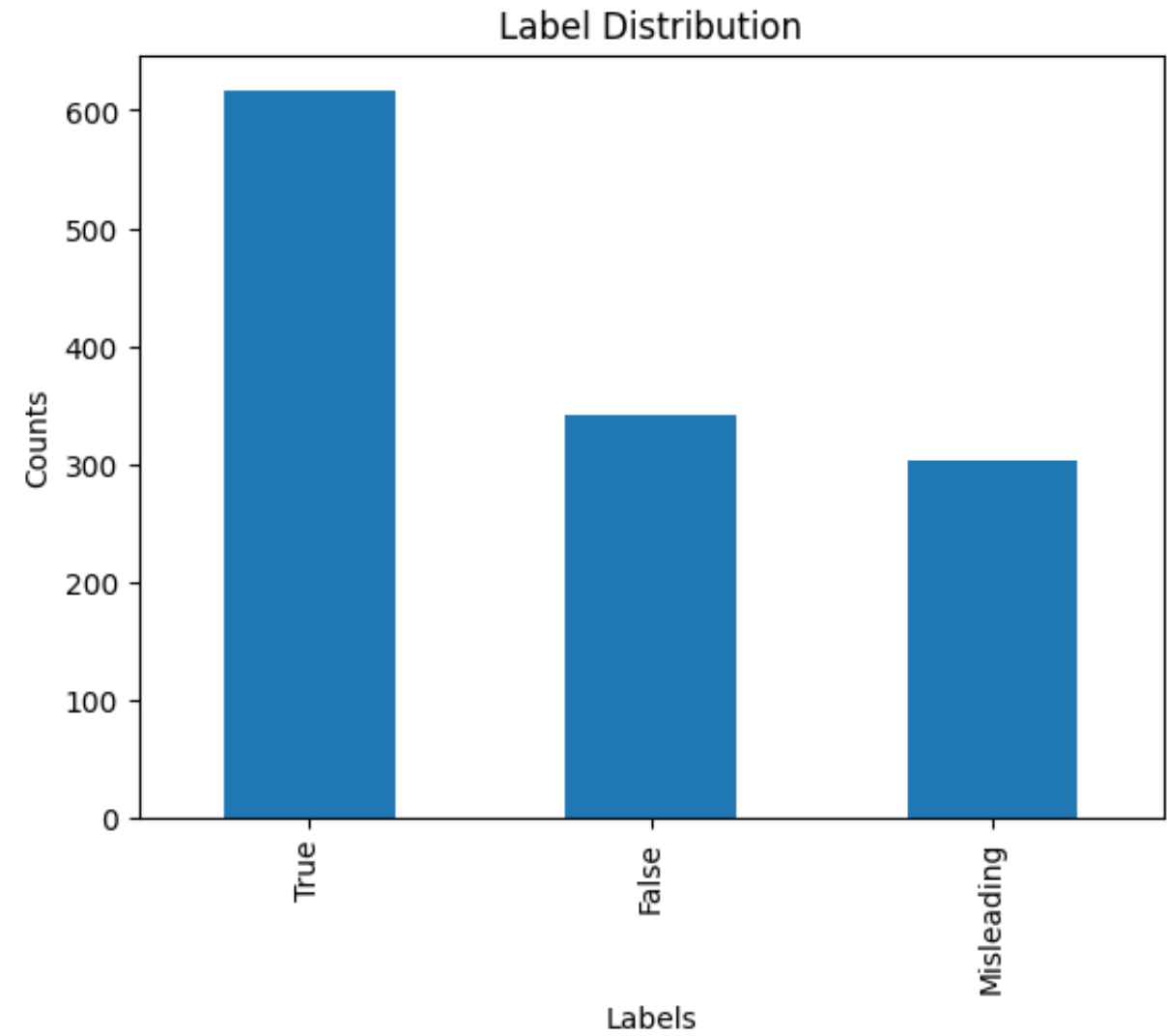
## Label Distribution



----------------------------------------

## 3. Training Strategy

```python
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='record
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', l
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. Training and Evaluation Before fine tuning using Gemini 2.0 Flash

A. Training and Evaluation Combination Between HealthFact and Scifact Dataset

```
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 train_claims = train_combined_data['claim'].tolist()
16 train_labels = train_combined_data['label'].tolist()
17 val_claims = val_combined_data['claim'].tolist()
```

```python
18 val_labels = val_combined_data['label'].tolist()
19
20 # Convert string labels to integers
21 label_encoder = LabelEncoder()
22 train_labels = label_encoder.fit_transform(train_labels)
23 val_labels = label_encoder.transform(val_labels)
24
25 # Load the BERT tokenizer
26 tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
27
28 # Tokenize the input data for pre-training
29 train_encodings = tokenizer(train_claims, truncation=True, padding=True, ma
30 val_encodings = tokenizer(val_claims, truncation=True, padding=True, max_le
31
32 # Create a dataset class
33 class ClaimsDataset(torch.utils.data.Dataset):
34     def __init__(self, encodings, labels):
35         self.encodings = encodings
36         self.labels = labels
37
38     def __getitem__(self, idx):
39         item = {key: torch.tensor(val[idx]) for key, val in self.encodings.
40         item['labels'] = torch.tensor(self.labels[idx])
41         return item
42
43     def __len__(self):
44         return len(self.labels)
45
46 # Create datasets for pre-training
47 train_dataset = ClaimsDataset(train_encodings, train_labels)
48 val_dataset = ClaimsDataset(val_encodings, val_labels)
49
50 # Load the BERT model
51 model = BertForSequenceClassification.from_pretrained('bert-base-uncased',
52
53 # Define training arguments for pre-training with validation loss logging
54 training_args = TrainingArguments(
55     output_dir='./results/pretrain',
56     num_train_epochs=3,
57     per_device_train_batch_size=8,
58     per_device_eval_batch_size=8,
59     warmup_steps=500,
60     weight_decay=0.01,
61     logging_dir='./logs/pretrain',
62     logging_steps=10,
63     eval_strategy="epoch",  # Updated to eval_strategy
64 )
65
66 # Create a Trainer instance for pre-training
67 trainer = Trainer(
68     model=model,
```

```
69       args=training_args,
70       train_dataset=train_dataset,
71       eval_dataset=val_dataset,
72       compute_metrics=lambda p: {
73           'accuracy': accuracy_score(p.label_ids, p.predictions.argmax(-1)),
74           'precision': precision_score(p.label_ids, p.predictions.argmax(-1),
75           'recall': recall_score(p.label_ids, p.predictions.argmax(-1), avera
76           'f1': f1_score(p.label_ids, p.predictions.argmax(-1), average='weig
77           'roc_auc': roc_auc_score(p.label_ids, torch.softmax(torch.tensor(p.
78       },
79 )
80
81 # Pre-train the model
82 trainer.train()
83
```

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 13:22, Epoch 3/3]

| Epoch | Training Loss | Validation Loss | Accuracy | Precision | Recall | F1 | Roc Auc |
|---|---|---|---|---|---|---|---|
| 1 | 0.945300 | 0.989038 | 0.547237 | 0.474881 | 0.547237 | 0.479945 | 0.886185 |
| 2 | 0.570800 | 0.998892 | 0.594177 | 0.571580 | 0.594177 | 0.557430 | 0.893968 |
| 3 | 0.443300 | 1.220652 | 0.590612 | 0.597319 | 0.590612 | 0.590918 | 0.896878 |

TrainOutput(global_step=4152, training_loss=0.7195591835990577, metrics=
{'train_runtime': 802.4998, 'train_samples_per_second': 41.364,

B. Training and Evaluation Covid Fake News Dataset

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

| Epoch | Training Loss | Validation Loss | Accuracy | Precision | Recall | F1 | Roc Auc |
|-------|---------------|-----------------|----------|-----------|--------|-----|---------|
| 1 | 0.425300 | 1.497399 | 0.566845 | 0.572513 | 0.566845 | 0.559901 | 0.874701 |
| 2 | 0.297900 | 1.842367 | 0.590018 | 0.596393 | 0.590018 | 0.591040 | 0.888562 |
| 3 | 0.208000 | 2.297663 | 0.588829 | 0.605546 | 0.588829 | 0.594764 | 0.889416 |

TrainOutput(global_step=4152, training_loss=0.3173761432093204, metrics=
{'train runtime': 809.0842, 'train samples per second': 41.028,

## ⌄ 5. Model Initialization

```
1 import google.generativeai as genai
2
3 # Initialize the Gemini 2.0 Flash Model:
4 API_KEY = "AIzaSyCjJIJ3ntglUAqZHn6qZaOE5uIGg4txEC4"
5 genai.configure(api_key=API_KEY)
6
7 # Load the Gemini model
8 model = genai.GenerativeModel("gemini-2.0-flash")
```

```
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 governme
9 """
10
11 response = generate_response(prompt)
12 print(response)
```

This claim appears to be a tweet or social media post promoting a segment o

**Here's a breakdown of the elements:**

* **"6 10"**: This likely refers to the time the tweet was posted, possibly
* **"Sky s"**:  This is likely a shortened form of "Sky News's".
* **"EdConwaySky"**: This is probably the Twitter handle for Ed Conway, who
* **"explains the latest COVID19 data and government announcement"**: This
* **"Get more on the coronavirus data here"**: This is a call to action, en
* **"https t co jvGZlSbFjH https t co PygSKXesBg"**: These are shortened UR

**Potential Issues and Things to Consider:**

* **Data Accuracy:**  While the claim itself isn't making a specific factua
* **Bias:**  It's important to be aware of potential biases. Sky News, like
* **Outdated Information:**  COVID-19 data and government announcements cha
* **Link Safety:** Always be cautious when clicking on shortened links, esp

**In Conclusion:**

The claim itself is a straightforward promotion of a news segment.  However

## ⌄ 5. Model Initialization

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

Requirement already satisfied: google-genai in /usr/local/lib/python3.11/di
Requirement already satisfied: anyio<5.0.0,>=4.8.0 in /usr/local/lib/python
Requirement already satisfied: google-auth<3.0.0,>=2.14.1 in /usr/local/lib
Requirement already satisfied: httpx<1.0.0,>=0.28.1 in /usr/local/lib/pytho
Requirement already satisfied: pydantic<3.0.0,>=2.0.0 in /usr/local/lib/pyt
Requirement already satisfied: requests<3.0.0,>=2.28.1 in /usr/local/lib/py
Requirement already satisfied: websockets<15.1.0,>=13.0.0 in /usr/local/lib
Requirement already satisfied: typing-extensions<5.0.0,>=4.11.0 in /usr/loc
Requirement already satisfied: idna>=2.8 in /usr/local/lib/python3.11/dist-
Requirement already satisfied: sniffio>=1.1 in /usr/local/lib/python3.11/di
Requirement already satisfied: cachetools<6.0,>=2.0.0 in /usr/local/lib/pyt
Requirement already satisfied: pyasn1-modules>=0.2.1 in /usr/local/lib/pyth
Requirement already satisfied: rsa<5,>=3.1.4 in /usr/local/lib/python3.11/d
Requirement already satisfied: certifi in /usr/local/lib/python3.11/dist-pa
Requirement already satisfied: httpcore==1.* in /usr/local/lib/python3.11/d
Requirement already satisfied: h11<0.15,>=0.13 in /usr/local/lib/python3.11
Requirement already satisfied: annotated-types>=0.6.0 in /usr/local/lib/pyt
Requirement already satisfied: pydantic-core==2.33.0 in /usr/local/lib/pyth
Requirement already satisfied: typing-inspection>=0.4.0 in /usr/local/lib/p
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/p
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3
Requirement already satisfied: pyasn1<0.7.0,>=0.6.1 in /usr/local/lib/pytho
Go to the following link in your browser, and complete the sign-in prompts:

    https://accounts.google.com/o/oauth2/auth?response_type=code&client_id=

Once finished, enter the verification code provided in your browser: 4/0Ab_

Credentials saved to file: [/content/.config/application_default_credential

These credentials will be used by any library that requests Application Def
WARNING:
Cannot find a quota project to add to ADC. You might receive a "quota excee

```
1 !pip install --upgrade google-cloud-aiplatform
```

Requirement already satisfied: google-cloud-aiplatform in /usr/local/lib/py
Requirement already satisfied: google-api-core!=2.0.*,!=2.1.*,!=2.2.*,!=2.3
Requirement already satisfied: google-auth<3.0.0,>=2.14.1 in /usr/local/lib
Requirement already satisfied: proto-plus<2.0.0,>=1.22.3 in /usr/local/lib/
Requirement already satisfied: protobuf!=4.21.0,!=4.21.1,!=4.21.2,!=4.21.3,
Requirement already satisfied: packaging>=14.3 in /usr/local/lib/python3.11
Requirement already satisfied: google-cloud-storage<3.0.0,>=1.32.0 in /usr/
Requirement already satisfied: google-cloud-bigquery!=3.20.0,<4.0.0,>=1.15.
Requirement already satisfied: google-cloud-resource-manager<3.0.0,>=1.3.3
Requirement already satisfied: shapely<3.0.0 in /usr/local/lib/python3.11/d
Requirement already satisfied: pydantic<3 in /usr/local/lib/python3.11/dist
Requirement already satisfied: typing-extensions in /usr/local/lib/python3.
Requirement already satisfied: docstring-parser<1 in /usr/local/lib/python3
Requirement already satisfied: googleapis-common-protos<2.0.0,>=1.56.2 in /
Requirement already satisfied: requests<3.0.0,>=2.18.0 in /usr/local/lib/py
Requirement already satisfied: grpcio<2.0dev,>=1.33.2 in /usr/local/lib/pyt
Requirement already satisfied: grpcio-status<2.0.dev0,>=1.33.2 in /usr/loca
Requirement already satisfied: cachetools<6.0,>=2.0.0 in /usr/local/lib/pyt
Requirement already satisfied: pyasn1-modules>=0.2.1 in /usr/local/lib/pyth
Requirement already satisfied: rsa<5,>=3.1.4 in /usr/local/lib/python3.11/d
Requirement already satisfied: google-cloud-core<3.0.0,>=2.4.1 in /usr/loca
Requirement already satisfied: google-resumable-media<3.0.0,>=2.0.0 in /usr
Requirement already satisfied: python-dateutil<3.0.0,>=2.8.2 in /usr/local/
Requirement already satisfied: grpc-google-iam-v1<1.0.0,>=0.14.0 in /usr/lo
Requirement already satisfied: google-crc32c<2.0dev,>=1.0 in /usr/local/lib
Requirement already satisfied: annotated-types>=0.6.0 in /usr/local/lib/pyt
Requirement already satisfied: pydantic-core==2.33.0 in /usr/local/lib/pyth
Requirement already satisfied: typing-inspection>=0.4.0 in /usr/local/lib/p
Requirement already satisfied: numpy<3,>=1.14 in /usr/local/lib/python3.11/
Requirement already satisfied: pyasn1<0.7.0,>=0.6.1 in /usr/local/lib/pytho
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.11/dist-p
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/p
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.11/di
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3

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

```
/usr/local/lib/python3.11/dist-packages/google/auth/_default.py:76: UserWar
    warnings.warn(_CLOUD_SDK_CREDENTIALS_WARNING)
INFO:vertexai.tuning._tuning:Creating SupervisedTuningJob
/usr/local/lib/python3.11/dist-packages/google/auth/_default.py:76: UserWar
    warnings.warn(_CLOUD_SDK_CREDENTIALS_WARNING)
INFO:vertexai.tuning._tuning:SupervisedTuningJob created. Resource name: pr
INFO:vertexai.tuning._tuning:To use this SupervisedTuningJob in another ses
INFO:vertexai.tuning._tuning:tuning_job = sft.SupervisedTuningJob('projects
INFO:vertexai.tuning._tuning:View Tuning Job:
https://console.cloud.google.com/vertex-ai/generative/language/locations/us
```

**⚏ VIEW TUNING JOB**

```
1 from google import genai
2 from google.genai import types
3 import base64
4
5 def generate():
6   client = genai.Client(
7       vertexai=True,
8       project="181085238689",
9       location="us-central1",
10   )
11
12   msg3_text1 = types.Part.from_text(text="""Clearly the Obama administratic
13
14   model = "projects/181085238689/locations/us-central1/endpoints/5419770989
15   contents = [
16     types.Content(
17       role="user",
18       parts=[
```

```python
19              types.Part.from_text(text="""Multiple Facebook posts claim that Aus
20          ]
21       ),
22       types.Content(
23          role="model",
24          parts=[
25             types.Part.from_text(text=label)
26          ]
27       ),
28       types.Content(
29          role="user",
30          parts=[
31             msg3_text1
32          ]
33       ),
34    ]
35    generate_content_config = types.GenerateContentConfig(
36       temperature = 0.2,
37       top_p = 0.8,
38       max_output_tokens = 1024,
39       response_modalities = ["TEXT"],
40       safety_settings = [types.SafetySetting(
41          category="HARM_CATEGORY_HATE_SPEECH",
42          threshold="OFF"
43       ),types.SafetySetting(
44          category="HARM_CATEGORY_DANGEROUS_CONTENT",
45          threshold="OFF"
46       ),types.SafetySetting(
47          category="HARM_CATEGORY_SEXUALLY_EXPLICIT",
48          threshold="OFF"
49       ),types.SafetySetting(
50          category="HARM_CATEGORY_HARASSMENT",
51          threshold="OFF"
52       )],
53    )
54
55    for chunk in client.models.generate_content_stream(
56       model = model,
57       contents = contents,
58       config = generate_content_config,
59       ):
60       print(chunk.text, end="")
61
62 generate()
```

```
/usr/local/lib/python3.11/dist-packages/google/auth/_default.py:76: UserWar
   warnings.warn(_CLOUD_SDK_CREDENTIALS_WARNING)
fake
```

```
1 Start coding or generate with AI.
```