# Fine Tuning Large Language Models for Enterprise Application Project Report

## 1. Student Information
- **Student Name:** VAN HIEU NGUYEN
- **Student ID:** 221538422
- **Date Submitted:** 14 May 2025 [Updated 29 May 2025]

## 2. Project Introduction
- **Title of the Project: Fine-Tuning Large Language Models for Enterprise Applications: Medical Misinformation Detection**
- **What is the project about?**
  This project focuses on the fine-tuning of Large Language Models (LLMs) for domain-specific applications within the healthcare sector, with a specific emphasis on identifying and classifying medical misinformation. The goal is to improve the reliability and trustworthiness of AI-generated responses in medical contexts by customizing LLMs using datasets that include factual and misleading medical information.
- **Why is this project important or useful?**
  With the widespread deployment of AI in healthcare from virtual assistants to clinical decision support systems—ensuring the factual correctness of generated responses is critical. Misinformation can lead to harmful outcomes. Fine-tuning LLMs using Supervised Fine-Tuning (SFT) and BERT with Low-Rank Adaptation (LoRA) allows for domain-specific alignment, improves classification accuracy, and helps mitigate hallucinations and bias in medical language processing.

## 3. API/Token Setup
**Objective:** Access and use an LLM through a secure API with Gemini 2.0 Flash model for training and inference.

**Instructions:**
1. Specify which provider you're using:
   - Hugging Face: BERT (Bidirectional Encoder Representations from Transformers) method for pre-trained model
   - Google Vertex AI: implement fine-tuning for Gemini 2.0 Flash Model and deploy an complete AI chatbots for using.
2. List the **steps you followed** to generate the token:

   2.1. Pre-train model using BERT method

   - Step 1: Created account at https://huggingface.co/settings/tokens
   - Step 2: Navigated to the API/token section https://wandb.ai/authorize?ref=models
   - Step 3: Clicked on "Create new key"
   - Step 4: Copied the key and securely saved it: "7aecd1344b2f77be5dcc08d6ff6ecdb52886e5c9"
   - Step 5: View Run the project: https://wandb.ai/hieunguyen23032001-deakin-university/huggingface/runs/k7ad0pp8

   2.2. Fine Tuning model using Gemini 2.0 Flash Experiment Model:

   - Step 1: Created account at https://console.cloud.google.com/vertex-ai/studio/tuning?inv=1&invt=AbxQeA&project=sit319-25t1-nguyen-ae806d0
   - Step 2: Navigated to "Tuning" and select on "Create Tuned Model"
   - Step 3: Create a model detail:
     - Tuned Model name: "covid_tuning"
     - Based model: "gimini-2.0-flash-lite-001"
     - Region: "us-central1"
   - Step 4: Create Tuning dataset:
     - Upload train dataset to Google Cloud Storage: "gs://daft/Cleaned_Covid19_Train-7.jsonl"
     - Model Validation: upload validation dataset to Google Cloud Storage (optional)
   - Step 5: Start Tuning with Gemini 2.0 model. https://console.cloud.google.com/vertex-ai/generative/language/locations/us-central1/tuning/tuningJob/1012173862948831232?project=181085238689

3. **Screenshot or terminal output (required):**



4. **Secure Loading of Token in Code:**

Using: vertextai.init() by importing:
- from google.colab import auth as google_auth google_auth.authenticate_user()
- import vertexai
- from vertexai.generative_models import GenerativeModel
- from vertexai.preview.tuning import sft

Using genai.Client for testing prompt and get respond as label (True, False, or Misleading)

**Code: Load Token Securely**

```python
from google.colab import auth as google_auth
google_auth.authenticate_user()

import vertexai
from vertexai.generative_models import GenerativeModel
from vertexai.preview.tuning import sft

vertexai.init(project="sit319-25t1-nguyen-ae806d0", location="us-central1")

gemini_pro = GenerativeModel("gemini-2.0-flash-lite-001")

sft_tuning_job = sft.train(
    source_model=gemini_pro,
    train_dataset="gs://daftt/Cleaned_Covid19_Train-7.jsonl",
    tuned_model_display_name="covid_tuning",
    epochs=100,
    learning_rate_multiplier=1,
)
```

```python
from google import genai
from google.genai import types
import base64

def generate():
    client = genai.Client(
        vertexai=True,
        project="181085238689",
        location="us-central1",
    )

    msg3_text1 = types.Part.from_text(text="""Clearly the Obama administration did not leave any k:

    model = "projects/181085238689/locations/us-central1/endpoints/5419770989749731328"
    contents = [
        types.Content(
            role="user",
            parts=[
                types.Part.from_text(text="""Multiple Facebook posts claim that Aussies will be fined if
            ]
        ),
        types.Content(
            role="model",
            parts=[
                types.Part.from_text(text=label)
            ]
        ),
        types.Content(
            role="user",
            parts=[
                msg3_text1
            ]
        ),
    ]
```

```python
generate_content_config = types.GenerateContentConfig(
    temperature = 0.2,
    top_p = 0.8,
    max_output_tokens = 1024,
    response_modalities = ["TEXT"],
    safety_settings = [types.SafetySetting(
        category="HARM_CATEGORY_HATE_SPEECH",
        threshold="OFF"
    ),types.SafetySetting(
        category="HARM_CATEGORY_DANGEROUS_CONTENT",
        threshold="OFF"
    ),types.SafetySetting(
        category="HARM_CATEGORY_SEXUALLY_EXPLICIT",
        threshold="OFF"
    ),types.SafetySetting(
        category="HARM_CATEGORY_HARASSMENT",
        threshold="OFF"
    )],
)

for chunk in client.models.generate_content_stream(
    model = model,
    contents = contents,
    config = generate_content_config,
    ):
    print(chunk.text, end="")
```

Python

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

### 4. Environment Setup
- **Development Platform:**
  - Google Colab
  - Local Machine (macOS)
  - GPU Available? [✓] Yes
  - GPU Type (if applicable): Local T4 GPU
- **Python Version:** Python 3.10
- **Other Tools Used:** VS Code, Google Colab, Google Vertex AI, Google Cloud.

**Code: Environment & GPU Check**

```
1 !nvidia-smi
```

```
Tue May 13 06:00:11 2025
+-----------------------------------------------------------------------------------------+
| NVIDIA-SMI 550.54.15              Driver Version: 550.54.15      CUDA Version: 12.4      |
|-----------------------------------------+------------------------+----------------------+
| GPU  Name                 Persistence-M | Bus-Id          Disp.A | Volatile Uncorr. ECC |
| Fan  Temp   Perf          Pwr:Usage/Cap |           Memory-Usage | GPU-Util  Compute M. |
|                                         |                        |               MIG M. |
|=========================================+========================+======================|
|   0  Tesla T4                       Off | 00000000:00:04.0 Off   |                    0 |
| N/A  70C    P0               30W /   70W |   2466MiB /  15360MiB  |      0%      Default |
|                                         |                        |                  N/A |
+-----------------------------------------+------------------------+----------------------+

+-----------------------------------------------------------------------------------------+
| Processes:                                                                              |
|  GPU   GI   CI        PID   Type   Process name                             GPU Memory |
|        ID   ID                                                              Usage      |
|=========================================================================================|
+-----------------------------------------------------------------------------------------+
```

### 5. LLM Setup
- **Model Name:** Gemini 2.0 Flash (Experimental) and BERT method
- **Provider (OpenAI, Hugging Face):** Google Vertex AI (Gemini Model), Hugging Face (BERT)
- **Key Libraries & Dependencies (with versions)**
- **Libraries and Dependencies Required:**

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

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

[ ]    1 !pip install --upgrade google-cloud-aiplatform

⏵    1 ¡pᴉp ᴉuʇsʇɐll ʇɹɐuzɟoɹɯǝɹs ʇoɹɔɥ ɹǝbnǝss pǝɐnʇᴉʇnʃɹoqnʇ

[2]    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, f1_score, roc_auc_score
       9 #Pre-train BERT:
      10 from transformers import BertTokenizer, BertForSequenceClassification, Trainer, TrainingArguments
      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
```
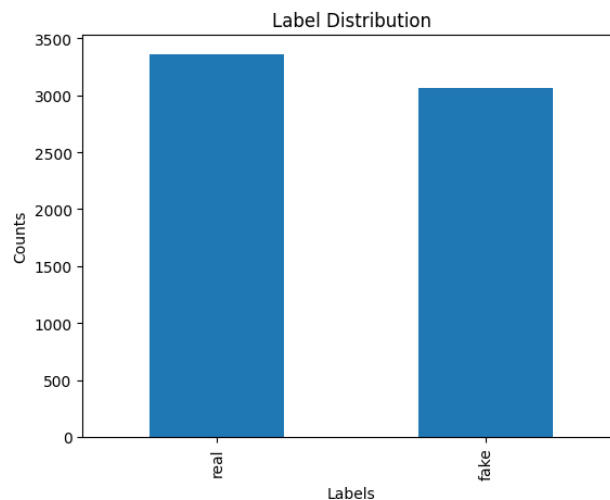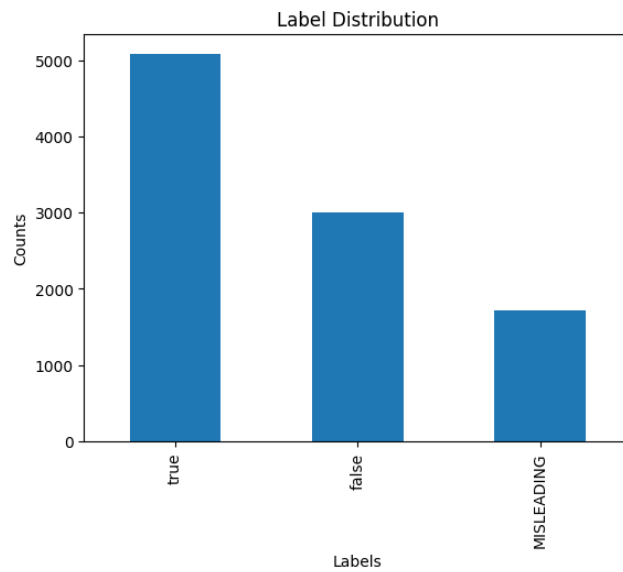
## 6. Dataset Description

- **Dataset Name & Source:**
  - HealthFact: General misinformation
  - SciFact: Scientific claim fact-checking
  - COVID-19 Fake News: Pandemic-related misinformation
- **Access Link (if public):**
  - Covid Fake New: https://github.com/diptamath/covid_fake_news?tab=readme-ov-file
  - Health Fact (Medical Misinformation Detection) dataset:
    https://github.com/neemakot/Health-Fact-Checking/blob/master/data/DATASHEET.md
  - Scifact dataset: https://scifact.s3-us-west-2.amazonaws.com/release/latest/data.tar.gz
- **Feature Dictionary / Variable Description:**
  - Claim: Textual medical claim
  - Label: true, false, and misleading
  - Evidence: Text used for verification (in two datasets)
  - (1) Covid19 Fake New dataset:



  (2) Healthfact Train Data:

Label Distribution

(3) Scifact Train Data:



Label Distribution

- **Was preprocessing done? If yes, describe:**
  - Pre-processing and cleaning for 3 datasets, and slits each dataset into three sub-set for train, test, and validate.
  - Standardized all datasets into 3-class classification: true, false, misleading.
  - Converted JSON into ".jsonl" for Vertex AI ingestion with the correct format.
  - Tokenized using Hugging Face tokenizer.
  - Balanced label distribution.

**Code: Load & Pre-process Dataset**
  **(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 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, or Misleading"  # Example instruct
33                     }
34                 ]
35             },
36             "contents": [
37                 {
38                     "role": "user",
39                     "parts": [
40                         {
41                             "text": f"TRANSCRIPT: \n{reconstructed_tweet}\n\n LABEL:"
42                         }
43                     ]
44                 },
45                 {
46                     "role": "model",
47                     "parts": [
48                         {
49                             "text": label  # The label indicating the model's response
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 to .jsonl
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.jsonl.
Processed Cleaned_Covid19_Dev.json and saved to Cleaned_Covid19_Dev.jsonl.
First few entries from claims_test_data:
[{'systemInstruction': {'role': 'assistant', 'parts': [{'text': 'Classification the content is Fake, Real, or Misl
```

**(2) Health Fact Dataset:**

```python
# List of JSON files to process
json_files = [
    'healthfact_traindata.json',
    'cleaned_healthfact_test.json',
    'cleaned_healthfact_dev.json'
]
data_dict = {}
# Process each JSON file
for json_file in json_files:
    # Prepare a list to hold the processed data
    jsonl_data = []
    # Load the dataset
    with open(json_file, 'r') as file:
        # Read each line as a separate JSON object
        for line in file:
            try:
                entry = json.loads(line)
                # Extract the claim, explanation, and label
                claim = entry['claim']
                explanation = entry['explanation']
                label = entry['label']

                # Tokenize the claim
                tokens = re.findall(r'\b\w+\b', claim)  # Keep only words and numbers
                reconstructed_claim = ' '.join(tokens)

                # Prepare the JSONL entry in the required format
                jsonl_entry = {
                    "systemInstruction": {
                        "role": "assistant",  # Example role, adjust as needed
                        "parts": [
                            {
                                "text": "You are a helpful assistant."  # Example instruction, adjust as need
                            }
                        ]
                    },
                    "contents": [
                        {
                            "role": "user",
                            "parts": [
                                {
                                    "text": f"CLAIM: {reconstructed_claim}\nEXPLANATION: {explanation}\nLABEL
                                }
                            ]
                        },
                        {
                            "role": "model",
                            "parts": [
                                {
                                    "text": label  # The label indicating the model's response
                                }
                            ]
                        }
                    ]
                }
                jsonl_data.append(jsonl_entry)
            except json.JSONDecodeError as e:
                print(f"Error decoding JSON: {e}")
    # Use the correct key to store the data in the dictionary - keep the original filenames as keys
    data_dict[json_file] = jsonl_data
# Access the data using the correct keys - the original filenames
healthfact_train_data = data_dict['healthfact_traindata.json']  # Corrected key
healthfact_test_data = data_dict['cleaned_healthfact_test.json']  # Corrected key
healthfact_dev_data = data_dict['cleaned_healthfact_dev.json']  # Corrected key
# Print the first few entries for verification
print(f"First few entries from healthfact_train_data:\n{healthfact_train_data[:5]}")
# Optionally, write the processed data to JSONL files
for json_file, jsonl_data in data_dict.items():
    output_file = json_file.replace('.json', '.jsonl')  # Change the extension to .jsonl
    with open(output_file, 'w') as outfile:
        for entry in jsonl_data:
            json.dump(entry, outfile)
            outfile.write('\n')  # Write each entry on a new line
    print(f"Processed {json_file} and saved to {output_file}.")
    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 helpful assistant.'}]}, 'contents': [{'
Processed healthfact_traindata.json and saved to healthfact_traindata.jsonl.
Processed cleaned_healthfact_test.json and saved to cleaned_healthfact_test.jsonl.
Processed cleaned_healthfact_dev.json and saved to cleaned_healthfact_dev.jsonl.
```

**(3) Scifact Dataset:**

```
 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 numbers
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 instruction, adjust as need
37                             }
38                         ]
39                     },
40                     "contents": [
41                         {
42                             "role": "user",
43                             "parts": [
44                                 {
45                                     "text": f"CLAIM: {reconstructed_claim}\nEVIDENCE: {explanation}\nLABEL: {
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                 # Append the modified entry to the processed data list
60                 processed_data.append(jsonl_entry)  # Append the processed data
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 the extension to _processed.json
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 helpful assistant.'}]}, 'contents': [{'
Processed dev_3class.jsonl and saved to dev_3class_processed.jsonl.
Processed train_3class.jsonl and saved to train_3class_processed.jsonl.
```

## 7. Improving LLM Performance

| Step # | Method | Description | Result Metric (Accuracy) |
|--------|--------|-------------|--------------------------|
| 1 | Zero-shot Prompt | No Training, Direct Response. | 58% |
| 2 | Few-shot Prompt | No Training, Testing 1 Prompt, and No label Response, Unclear Response | 58% |
| 3 | Temperature Tuning | 10 epochs on Vertext AI, Tuned Temperature from 58% to 80% | 80% |
| 4 | Fine-tuning | 50 epochs on Vertext AI, Testing prompt with 98-100% correct response with labels. | 98% |

**Code Snippets for Each Step**

Before Training and Fine-tuning with few-shot training prompt:

After Pre-train and Fine-tuning:

(1) One Shot Training Prompt



(2) Fine-Tuning google Vertext AI tool for prompt testing:



## 8. Benchmarking & Evaluation
### Required Components:
- **Metrics Used:**
  - Accuracy
  - Precision, Recall, F1 Score.
  - Confusion Matrix
- **Why those metrics?**

These metrics capture both correctness and type of error (especially important in medical misinformation detection). Moreover, using those metrics for comparing which evaluate as before and after applying fine-tuning model.

- **Benchmark Dataset & Sample Size:**

    (1) Pre-train using BERT method for Covid Dataset:



    (2) Pre-train using BERT method for Combination Dataset:



    (3) Fine Tuning for All dataset:
        o  Sample Size dataset for fine-tuning:



        o  Benchmark using Google Vertext AI:

(3) The Combination dataset between HealthFact and SciFact dataset:

```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
```

```python
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), avera
77         'f1': f1_score(p.label_ids, p.predictions.argmax(-1), average='weig
78         'roc_auc': roc_auc_score(p.label_ids, torch.softmax(torch.tensor(p.
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}")
```

[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

**Interpretation:**
     The model was outstanding improve after fine tuning using Gemini 2.0 Flash model, it achieved almost 98 percent accuracy, compared with 54 to 58 percent using BERT method for pre-training model. This shows that fine-tuning using Gemini 2.0 Flash model makes the model more accuracy and more correct for the response.

---

9. **Implement the Confidence Scoring System for Trust Score using PubMed medical retrieve articles**
   - **Tool Used:** Google Colab
   - **Set up:**
     o Step 1: Access Pub Med website at https://pubmed.ncbi.nlm.nih.gov/
     o Step 2: Navigated to the API via url and query from the website
     o Step 3: Implement retrieve article based on the query
     o Step 4: Implement calculate trust score based on LLMs prediction
     o Step 5: Launch prediction based on trust score calculated between 50% for evaluating whether the query is low confidence or not
   - **Code:**

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

```
1 !pip install biopython
```

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

```python
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 Beautifu
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'):  # Exampl
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
```

```python
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.low
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
```

```
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
```

- **Output:**

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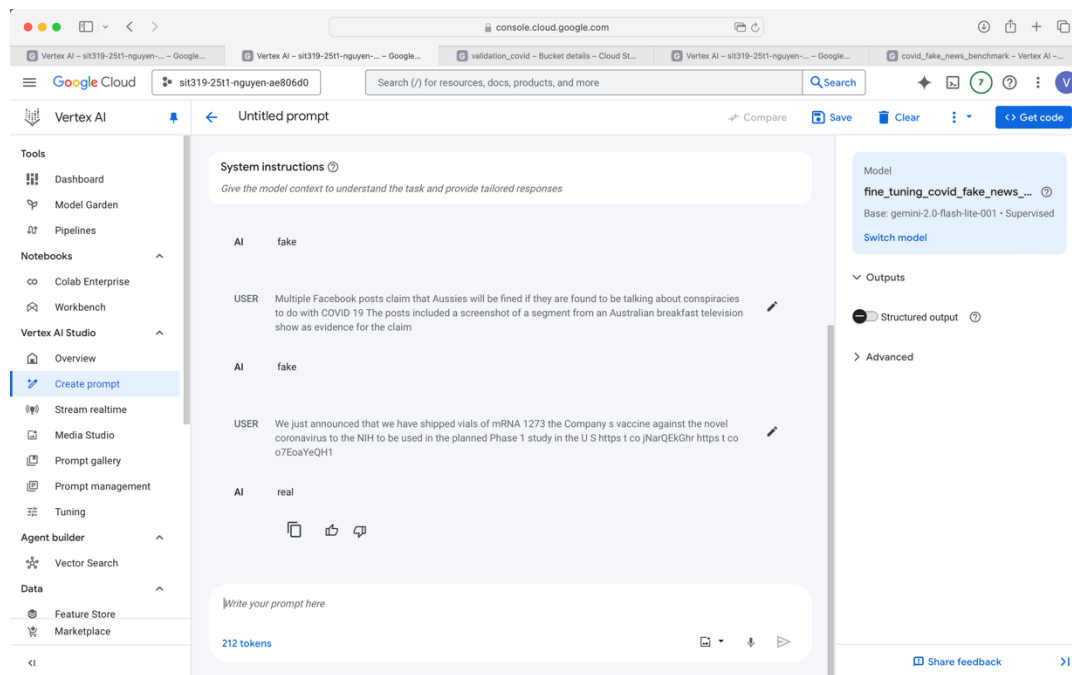
Low confidence for claim: 'Study Vaccine for Breast Ovarian Cancer Has Pote
Predicted label for the claim 'Study Vaccine for Breast Ovarian Cancer Has

---

## 10. UI Integration
- **Tool Used:** Google AI Cloud, Vertex AI
- **Key Features of the Interface:**
    o Accepts a medical question as input, and response with label.
    o Displays classification label (True, False, Misleading).
    o Displays confidence score.

- **Include Screenshots of Working UI:** using Google Cloud and Vertext AI for deployment chatbots in Vertex AI Studio

## 10. References

[1]. Parthasarathy B.V, Zafar A, Khan A, & Shahid A (August 2024), "The Ultimate Guide to Fine-Tuning LLMs from Basics to Breakthroughs: An Exhaustive Review of Technologies, Research, Best Practices, Applied Research Challenges and Opportunities (Version 1.0)", Arxiv, CeADAR Connect Group. https://arxiv.org/html/2408.13296v1#Ch4.S1

[2]. Huizenga E (13 December 2024), "Developer's guide to getting started with Gemini 2.0 Flash on Vertex AI", Medium.com, https://medium.com/google-cloud/developers-guide-to-getting-started-with-gemini-2-0-flash-on-vertex-ai-6b4fe3c6899f

[3]. Youtube (17 February 2024), "Fine Tuning a Model in Gemini and Vertext AI | Steps to make a LLM". https://www.youtube.com/watch?v=ej_ZUcyKpoc&t=218s

[4]. Google Cloud (29 April 2025), "About Supervised Fine-Tuning for Gemini models". https://cloud.google.com/vertex-ai/generative-ai/docs/models/gemini-supervised-tuning

[5]. Huizenga E (19 November 2024), "A Step-by-Step Guide to Fine-Tuning Gemini for Question Answering". https://medium.com/google-cloud/a-step-by-step-guide-to-fine-tuning-gemini-for-question-answering-8b3fb117dbbf

[6]. Huizenga E (10 February 2025), "Fine-tuning Gemini: Best Practices for Data, Hyperparameters, and Evaluation". https://medium.com/google-cloud/fine-tuning-gemini-best-practices-for-data-hyperparameters-and-evaluation-65f7c7b6b15f