

```

1 # IMPORTANT: SOME KAGGLE DATA SOURCES ARE PRIVATE
2 # RUN THIS CELL IN ORDER TO IMPORT YOUR KAGGLE DATA SOURCES.
3 import kagglehub
4 kagglehub.login()
5
6
7 # IMPORTANT: RUN THIS CELL IN ORDER TO IMPORT YOUR KAGGLE DATA SOURCES,
8 # THEN FEEL FREE TO DELETE THIS CELL.
9 # NOTE: THIS NOTEBOOK ENVIRONMENT DIFFERS FROM KAGGLE'S PYTHON
10 # ENVIRONMENT SO THERE MAY BE MISSING LIBRARIES USED BY YOUR
11 # NOTEBOOK.
12
13 aditya0506_dataset_zip_path = kagglehub.dataset_download('aditya0506/dataset-zip')
14 metaresearch_llama_3_2_transformers_3b_instruct_1_path = kagglehub.model_download('metaresearch/llama-3.2/Transformers/3b-instruct-1')
15 print('Data source import complete.')
16

```

```

1 %capture
2 %pip install -U bitsandbytes
3 %pip install -U transformers
4 %pip install -U accelerate
5 %pip install -U peft
6 %pip install -U trl

```

```

1 import numpy as np
2 import pandas as pd
3 import re
4 import os
5 from tqdm import tqdm
6 import bitsandbytes as bnb
7 import torch
8 import torch.nn as nn
9 import transformers
10 from datasets import Dataset
11 from peft import LoraConfig, PeftConfig, PeftModel
12 from trl import SFTTrainer
13 from trl import setup_chat_format
14 from transformers import (AutoModelForCausalLM,
15                           AutoTokenizer,
16                           BitsAndBytesConfig,
17                           TrainingArguments,
18                           pipeline,
19                           logging)
20 from sklearn.metrics import (accuracy_score,
21                             classification_report,
22                             confusion_matrix)
23 from sklearn.model_selection import train_test_split

```

2025-04-16 09:25:23.367876: E external/local\_xla/xla/stream\_executor/cuda/cuda\_fft.cc:477] Unable to register cuFFT factory: Attempting to register factory  
 WARNING: All log messages before absl::InitializeLog() is called are written to STDERR  
 E0000 00:00:1744795523.541222 31 cuda\_dnn.cc:8310] Unable to register cuDNN factory: Attempting to register factory  
 E0000 00:00:1744795523.597096 31 cuda\_blas.cc:1418] Unable to register cuBLAS factory: Attempting to register factory

```

1 #Preparing the dataset
2 data_dir = '/kaggle/input/dataset-zip/Dataset'
3 data_records = []
4
5 def normalize_citations_consistent(text):
6     """
7     Normalize citations by converting different citation formats
8     into a standardized placeholder "[Citation]".
9     Detects numeric citations like "[12]" or "[12, 34]" and
10    author-year citations like "(Smith et al., 2020)" or "Smith et al. (2020)".
11    """
12    # Replace numeric citations in square brackets.
13    text = re.sub(r'(\s*\d+(?:\s*,\s*\d+)*\s*)', '[Citation]', text)
14    # Replace citations in the form "(Smith et al., 2020)" or "Smith, 2020)".
15    text = re.sub(r'([A-Z][a-zA-Z]+(?: et al\.)?\s*\d{4})', '[Citation]', text)
16    # Also handle cases like "Smith et al. (2020)" without outer parentheses.
17    text = re.sub(r'[A-Z][a-zA-Z]+(?: et al\.)?\s+\s*\d{4}\s*', '[Citation]', text)
18    # Normalize excessive whitespace.
19    text = re.sub(r'\s+', ' ', text).strip()
20    return text
21
22 def extract_title_and_abstract(text):
23     """
24     Extracts the title and abstract from text using 'Title:' and 'Abstract:' keywords.

```

```

25 Falls back to treating the whole text as abstract if title is not found.
26
27 Returns:
28 (title, abstract): Tuple of extracted title and abstract
29 """
30 title = ""
31 abstract = text.strip()
32
33 # Normalize line endings
34 text = text.replace('\r\n', '\n').replace('\r', '\n')
35
36 # Extract title (case-insensitive, from first "Title:")
37 title_match = re.search(r'\bTitle:\s*(.*)', text, re.IGNORECASE)
38 if title_match:
39     title = title_match.group(1).strip()
40     # Remove title line from text
41     text = re.sub(r'\bTitle:\s*.*\n?', '', text, flags=re.IGNORECASE, count=1).strip()
42
43 # Extract abstract (case-insensitive, from first "Abstract:")
44 abstract_match = re.search(r'\bAbstract:\s*(.*)', text, re.IGNORECASE | re.DOTALL)
45 if abstract_match:
46     abstract = abstract_match.group(1).strip()
47
48 return title, abstract
49
50 # Iterate through each subfolder (e.g. folders named "cancer" or "non-cancer").
51 for folder in os.listdir(data_dir):
52     folder_path = os.path.join(data_dir, folder)
53     if os.path.isdir(folder_path):
54         # Determine the label based on folder name.
55         label = "Noncancer" if "non-cancer" in folder.lower() else "Cancer"
56         for file in os.listdir(folder_path):
57             if file.endswith('.txt'):
58                 file_path = os.path.join(folder_path, file)
59                 with open(file_path, 'r', encoding='utf-8') as f:
60                     text = f.read().strip()
61                     # Normalize citations.
62                     text = normalize_citations_consistent(text)
63                     # Extract title and abstract.
64                     title, abstract_text = extract_title_and_abstract(text)
65                     abstract_id = os.path.splitext(file)[0]
66                     data_records.append({
67                         "text": re.sub(r"<ID:[^>]+>\s*", "", abstract_text),
68                         "label": label
69                     })
70
71 # Create a Pandas DataFrame with the loaded data:
72 df = pd.DataFrame(data_records)
73 df.head()

```



	text	label
0	Title: Clear Cell Variant of Papillary Thyroid...	Cancer
1	Title: [Hereditary Breast and Ovarian Cancer S...	Cancer
2	Title: [Hepatic epithelioid angiomyolipoma/PEC...	Cancer
3	Title: Giant Lumbar Polypoid Tumor with Bullae...	Cancer
4	Title: Multidisciplinary Challenges in Mastocy...	Cancer

```

1 # Shuffle the DataFrame
2 df = df.sample(frac=1, random_state=85).reset_index(drop=True)
3
4 # Split the DataFrame
5 train_size = 0.8
6 eval_size = 0.1
7
8 # Calculate sizes
9 train_end = int(train_size * len(df))
10 eval_end = train_end + int(eval_size * len(df))
11
12 # Split the data
13 X_train = df[:train_end]
14 X_eval = df[train_end:eval_end]
15 X_test = df[eval_end:]
16
17 # Define the prompt generation functions
18 def generate_prompt(data_point):
19     return f"""
20         Classify the text into Cancer, Noncancer and return the answer as the corresponding research paper label.
21 text: {data_point["text"]}

```

```

22 label: {data_point["label"]}"".strip()
23
24 def generate_test_prompt(data_point):
25     return f"""
26         Classify the text into Cancer, Noncancer and return the answer as the corresponding research paper label.
27 text: {data_point["text"]}
28 label: "".strip()
29
30 # Generate prompts for training and evaluation data
31 X_train.loc[:, 'text'] = X_train.apply(generate_prompt, axis=1)
32 X_eval.loc[:, 'text'] = X_eval.apply(generate_prompt, axis=1)
33
34 # Generate test prompts and extract true labels
35 y_true = X_test.loc[:, 'label']
36 X_test = pd.DataFrame(X_test.apply(generate_test_prompt, axis=1), columns=["text"])

```

```
1 X_train.label.value_counts()
```

```

↗ label
  Cancer      410
 Noncancer    390
  Name: count, dtype: int64

```

```
1 y_true.value_counts()
```

```

↗ label
 Noncancer     55
   Cancer      45
  Name: count, dtype: int64

```

```

1 # Convert to datasets
2 train_data = Dataset.from_pandas(X_train[["text"]])
3 eval_data = Dataset.from_pandas(X_eval[["text"]])

```

```
1 train_data['text'][3]
```

```

↗ "Classify the text into Cancer, Noncancer and return the answer as the corresponding research paper label.\n
text: Title: Neurological disease in xeroderma pigmentosum: prospective cohort study of its features and progression.
Abstract: Xeroderma pigmentosum (XP) results from biallelic mutations in any of eight genes involved in DNA repair
systems, thus defining eight different genotypes (XPA, XPB, XPC, XPD, XPE, XPF, XPG and XP variant or XPV). In addition
to cutaneous and ophthalmological features, some patients present with XP neurological disease. It is unknown whether
the different neurological signs and their progression differ among groups. Therefore, we aim to characterize the XP
neurological disease and its evolution in the heterogeneous UK XP cohort. Patients with XP were followed in the UK
National XP Service, from 2009 to 2021. Age of onset for different events was recorded. Cerebellar ataxia and
additional neurological signs and symptoms were rated with the Scale for the Assessment and Rating of Ataxia (SARA),
the Inventory of Non-Ataxia Signs (INAS) and the Activities of Daily Living questionnaire (ADL). Patients' mutations
received scores based on their predicted effects. Data from available ancillary tests were collected. Ninety-three XP
patients were recruited. Thirty-six (38.7%) reported neurological symptoms, especially in the XPA, XPD and XPG groups,
with early-onset and late-onset forms, and typically appearing after cutaneous and ophthalmological symptoms. XPA, XPD
and XPG patients showed higher SARA scores compared to XPC, XPE and XPV. SARA total scores significantly increased over
time in XPD (0.91 points/year, 95% confidence interval: 0.61, 1.21) and XPA (0.63 points/year, 95% confidence interval:
0.38, 0.89). Hyporeflexia, hypopallesthesia, upper motor neuron signs, chorea, dystonia, oculomotor signs and
cognitive impairment were frequent findings in XPA, XPD and XPG. Cerebellar and global brain atrophy, axonal sensory
and sensorimotor neuropathies, and sensorineural hearing loss were common findings in patients. Some XPC, XPE and XPV
cases presented with abnormalities on examination and/or ancillary tests, suggesting underlying neurological
involvement. More severe mutations were associated with a faster progression in SARA total score in XPA (0.40
points/year per 1-unit increase in severity score) and XPD (0.60 points/year per 1-unit increase), and in ADL total
score in XPA (0.35 points/year per 1-unit increase). Symptomatic and asymptomatic forms of neurological disease are
frequent in XP patients, and neurological symptoms can be an important cause of disability. Typically, the neurological
disease will be preceded by cutaneous and ophthalmological features, and these should be actively searched in patients
with idiopathic late-onset neurological syndromes. Scales assessing cerebellar function, especially walking and speech,
and disability can show progression in some of the groups. Mutation severity can be used as a prognostic biomarker for
stratification purposes in clinical trials.\n
label: Noncancer"

```

## ✓ Loading the model and tokenizer

```

1 base_model_name = "/kaggle/input/llama-3.2/transformers/3b-instruct/1"
2
3 bnb_config = BitsAndBytesConfig(
4     load_in_4bit=True,
5     bnb_4bit_use_double_quant=False,
6     bnb_4bit_quant_type="nf4",
7     bnb_4bit_compute_dtype="float16",
8 )
9
10 model = AutoModelForCausalLM.from_pretrained(
11     base_model_name,
12     device_map="auto",
13     torch_dtype="float16",
14     quantization_config=bnb_config,

```

```

15 )
16
17 model.config.use_cache = False
18 model.config.pretraining_tp = 1

```

loading checkpoint shards: 0% | 0/2 [00:00<?. ?it/s]

```

1 tokenizer = AutoTokenizer.from_pretrained(base_model_name)
2
3 tokenizer.pad_token_id = tokenizer.eos_token_id

```

### Model evaluation before fine-tuning

```

1 match_labels = re.compile(r'\b(Non[-\s]?Cancer|Cancer)\b', flags=re.IGNORECASE)
2
3 def predict(test, model, tokenizer):
4     y_pred = []
5
6     for i in tqdm(range(len(test))):
7         prompt = test.iloc[i]["text"]
8         pipe = pipeline(task="text-generation",
9                         model=model,
10                        tokenizer=tokenizer,
11                        max_new_tokens=5,
12                        temperature=0.1)
13
14         result = pipe(prompt)
15         answer = result[0]['generated_text'].split("label: ")[-1].strip()
16         generated_text = result[0]['generated_text']
17         # Apply regex to extract the label
18         matches = match_labels.findall(answer)
19         if matches:
20             # Normalize and pick the first match
21             label = matches[0].replace("-", "").replace(" ", "").capitalize()
22             if label == "Noncancer":
23                 y_pred.append("Noncancer")
24             else:
25                 y_pred.append("Cancer")
26         else:
27             y_pred.append("none")
28
29     return y_pred

```

```

1 y_pred = predict(X_test, model, tokenizer)

```



```

81%|██████████| 81/100 [00:40<00:09, 1.92it/s]Device set to use cuda:0
82%|██████████| 82/100 [00:47<00:09, 1.93it/s]Device set to use cuda:0
83%|██████████| 83/100 [00:47<00:09, 1.78it/s]Device set to use cuda:0
84%|██████████| 84/100 [00:48<00:09, 1.70it/s]Device set to use cuda:0
85%|██████████| 85/100 [00:48<00:08, 1.77it/s]Device set to use cuda:0
86%|██████████| 86/100 [00:49<00:08, 1.72it/s]Device set to use cuda:0
87%|██████████| 87/100 [00:49<00:06, 1.90it/s]Device set to use cuda:0
88%|██████████| 88/100 [00:50<00:06, 1.73it/s]Device set to use cuda:0
89%|██████████| 89/100 [00:51<00:06, 1.72it/s]Device set to use cuda:0
90%|██████████| 90/100 [00:51<00:05, 1.67it/s]Device set to use cuda:0
91%|██████████| 91/100 [00:52<00:05, 1.73it/s]Device set to use cuda:0
92%|██████████| 92/100 [00:54<00:07, 1.09it/s]Device set to use cuda:0
93%|██████████| 93/100 [00:54<00:05, 1.20it/s]Device set to use cuda:0
94%|██████████| 94/100 [00:55<00:04, 1.25it/s]Device set to use cuda:0
95%|██████████| 95/100 [00:55<00:03, 1.36it/s]Device set to use cuda:0
96%|██████████| 96/100 [00:56<00:02, 1.48it/s]Device set to use cuda:0
97%|██████████| 97/100 [00:57<00:01, 1.53it/s]Device set to use cuda:0
98%|██████████| 98/100 [00:57<00:01, 1.45it/s]Device set to use cuda:0
99%|██████████| 99/100 [00:58<00:00, 1.50it/s]Device set to use cuda:0
100%|██████████| 100/100 [00:58<00:00, 1.70it/s]

```

```

1 def evaluate(y_true, y_pred):
2     labels = ["Noncancer", "Cancer"]
3     mapping = {label: idx for idx, label in enumerate(labels)}
4
5     def map_func(x):
6         return mapping.get(x, -1) # Map to -1 if not found, but should not occur with correct data
7
8     y_true_mapped = np.vectorize(map_func)(y_true)
9     y_pred_mapped = np.vectorize(map_func)(y_pred)
10
11     # Calculate accuracy
12     accuracy = accuracy_score(y_true=y_true_mapped, y_pred=y_pred_mapped)
13     print(f'Accuracy: {accuracy:.3f}')
14
15     # Generate accuracy report
16     unique_labels = set(y_true_mapped) # Get unique labels
17
18     for label in unique_labels:
19         label_indices = [i for i in range(len(y_true_mapped)) if y_true_mapped[i] == label]
20         label_y_true = [y_true_mapped[i] for i in label_indices]
21         label_y_pred = [y_pred_mapped[i] for i in label_indices]
22         label_accuracy = accuracy_score(label_y_true, label_y_pred)
23         print(f'Accuracy for label {labels[label]}: {label_accuracy:.3f}')
24
25     # Generate classification report
26     class_report = classification_report(y_true=y_true_mapped, y_pred=y_pred_mapped, target_names=labels, labels=list(range(len(labels))))
27     print('\nClassification Report:')
28     print(class_report)
29
30     # Generate confusion matrix
31     conf_matrix = confusion_matrix(y_true=y_true_mapped, y_pred=y_pred_mapped, labels=list(range(len(labels))))
32     print('\nConfusion Matrix:')
33     print(conf_matrix)

```

```

1 #Not FineTuned
2 evaluate(y_true, y_pred)

```

```

Accuracy: 0.390
Accuracy for label Noncancer: 0.255
Accuracy for label Cancer: 0.556

```

```

Classification Report:

```

	precision	recall	f1-score	support
Noncancer	0.45	0.25	0.33	55
Cancer	0.38	0.56	0.45	45
micro avg	0.40	0.39	0.40	100
macro avg	0.42	0.41	0.39	100
weighted avg	0.42	0.39	0.38	100

```

Confusion Matrix:
[[14 41]
 [17 25]]

```

```

1 import bitsandbytes as bnb
2 #Extract All the linear modules name
3 def find_all_linear_names(model):
4     cls = bnb.nn.Linear4bit
5     lora_module_names = set()
6     for name, module in model.named_modules():
7         if isinstance(module, cls):

```

```

8         names = name.split('.')
9         lora_module_names.add(names[0] if len(names) == 1 else names[-1])
10    if 'lm_head' in lora_module_names: # needed for 16 bit
11        lora_module_names.remove('lm_head')
12    return list(lora_module_names)

```

```

1 modules = find_all_linear_names(model)
2 modules

```

```

→ ['gate_proj', 'v_proj', 'o_proj', 'up_proj', 'down_proj', 'k_proj', 'q_proj']

```

```

1 #Model Set up
2 from trl import SFTTrainer, SFTConfig
3 output_dir="llama-3.2-fine-tuned-model"
4
5 peft_config = LoraConfig(
6     lora_alpha=16,
7     lora_dropout=0.1,
8     r=8,
9     bias="none",
10    task_type="CAUSAL_LM",
11    target_modules=modules,
12 )
13
14 training_arguments = SFTConfig(
15     output_dir="logs",
16     num_train_epochs=1,
17     gradient_checkpointing=True,
18     per_device_train_batch_size=1,
19     gradient_accumulation_steps=8,
20     optim="paged_adamw_32bit", # Use fused AdamW optimizer
21     save_steps=100,
22     load_best_model_at_end=True,
23     logging_steps=25,
24     learning_rate=2e-4,
25     weight_decay=0.001,
26     fp16=True,
27     bfloat16=False,
28     max_grad_norm=0.3,
29     max_steps=-1,
30     warmup_ratio=0.03,
31     group_by_length=False,
32     save_strategy="steps",
33     eval_steps=100,
34     eval_accumulation_steps=1,
35     lr_scheduler_type="cosine",
36     report_to="tensorboard",
37     eval_strategy="steps", # save checkpoint every epoch
38     max_seq_length=512,
39     packing=False,
40     dataset_kwargs={
41         "add_special_tokens": False, # Template with special tokens
42         "append_concat_token": False, # Add EOS token as separator token
43     }
44 )
45
46
47 trainer = SFTTrainer(
48     model=model,
49     train_dataset=train_data,
50     eval_dataset=eval_data,
51     peft_config=peft_config,
52     processing_class=tokenizer,
53     args=training_arguments,
54 )

```

```

→ Converting train dataset to ChatML: 0%|          | 0/800 [00:00<?, ? examples/s]
Applying chat template to train dataset: 0%|          | 0/800 [00:00<?, ? examples/s]
Tokenizing train dataset: 0%|          | 0/800 [00:00<?, ? examples/s]
Truncating train dataset: 0%|          | 0/800 [00:00<?, ? examples/s]
Converting eval dataset to ChatML: 0%|          | 0/100 [00:00<?, ? examples/s]
Applying chat template to eval dataset: 0%|          | 0/100 [00:00<?, ? examples/s]
Tokenizing eval dataset: 0%|          | 0/100 [00:00<?, ? examples/s]
Truncating eval dataset: 0%|          | 0/100 [00:00<?, ? examples/s]
No label_names provided for model class `PeftModelForCausalLM`. Since `PeftModel` hides base models input arguments, if

```

```

1 # Train model
2 trainer.train()

```



[100/100 21:28, Epoch 1/1]

Step Training Loss Validation Loss

100 1.833500 1.768290

```
TrainOutput(global_step=100, training_loss=1.867441635131836, metrics={'train_runtime': 1301.2147,
'train samples per second': 0.615, 'train steps per second': 0.077, 'total flos': 5103678867886080.0, 'train loss':
```

```
1 y_pred = predict(X_test, model, tokenizer)
2 evaluate(y_true, y_pred)
```



```
61%|██████████| 61/100 [00:39<00:28, 1.36it/s]Device set to use cuda:0
62%|██████████| 62/100 [00:40<00:26, 1.43it/s]Device set to use cuda:0
63%|██████████| 63/100 [00:41<00:26, 1.39it/s]Device set to use cuda:0
64%|██████████| 64/100 [00:41<00:25, 1.44it/s]Device set to use cuda:0
65%|██████████| 65/100 [00:42<00:23, 1.51it/s]Device set to use cuda:0
66%|██████████| 66/100 [00:43<00:24, 1.37it/s]Device set to use cuda:0
67%|██████████| 67/100 [00:43<00:25, 1.32it/s]Device set to use cuda:0
68%|██████████| 68/100 [00:44<00:22, 1.40it/s]Device set to use cuda:0
69%|██████████| 69/100 [00:45<00:21, 1.43it/s]Device set to use cuda:0
70%|██████████| 70/100 [00:45<00:20, 1.45it/s]Device set to use cuda:0
71%|██████████| 71/100 [00:46<00:19, 1.49it/s]Device set to use cuda:0
72%|██████████| 72/100 [00:47<00:18, 1.53it/s]Device set to use cuda:0
73%|██████████| 73/100 [00:47<00:18, 1.44it/s]Device set to use cuda:0
74%|██████████| 74/100 [00:48<00:20, 1.26it/s]Device set to use cuda:0
75%|██████████| 75/100 [00:49<00:19, 1.26it/s]Device set to use cuda:0
76%|██████████| 76/100 [00:50<00:17, 1.36it/s]Device set to use cuda:0
77%|██████████| 77/100 [00:51<00:17, 1.35it/s]Device set to use cuda:0
78%|██████████| 78/100 [00:51<00:14, 1.49it/s]Device set to use cuda:0
79%|██████████| 79/100 [00:52<00:15, 1.39it/s]Device set to use cuda:0
80%|██████████| 80/100 [00:52<00:12, 1.58it/s]Device set to use cuda:0
81%|██████████| 81/100 [00:53<00:11, 1.60it/s]Device set to use cuda:0
82%|██████████| 82/100 [00:53<00:10, 1.72it/s]Device set to use cuda:0
83%|██████████| 83/100 [00:54<00:10, 1.55it/s]Device set to use cuda:0
84%|██████████| 84/100 [00:55<00:11, 1.45it/s]Device set to use cuda:0
85%|██████████| 85/100 [00:56<00:09, 1.60it/s]Device set to use cuda:0
86%|██████████| 86/100 [00:56<00:09, 1.50it/s]Device set to use cuda:0
87%|██████████| 87/100 [00:57<00:08, 1.61it/s]Device set to use cuda:0
88%|██████████| 88/100 [00:57<00:07, 1.57it/s]Device set to use cuda:0
89%|██████████| 89/100 [00:58<00:06, 1.62it/s]Device set to use cuda:0
90%|██████████| 90/100 [00:59<00:06, 1.51it/s]Device set to use cuda:0
91%|██████████| 91/100 [00:59<00:05, 1.50it/s]Device set to use cuda:0
92%|██████████| 92/100 [01:01<00:08, 1.04s/it]Device set to use cuda:0
93%|██████████| 93/100 [01:02<00:06, 1.04it/s]Device set to use cuda:0
94%|██████████| 94/100 [01:03<00:05, 1.07it/s]Device set to use cuda:0
95%|██████████| 95/100 [01:04<00:04, 1.15it/s]Device set to use cuda:0
96%|██████████| 96/100 [01:04<00:03, 1.23it/s]Device set to use cuda:0
97%|██████████| 97/100 [01:05<00:02, 1.27it/s]Device set to use cuda:0
98%|██████████| 98/100 [01:06<00:01, 1.20it/s]Device set to use cuda:0
99%|██████████| 99/100 [01:07<00:00, 1.23it/s]Device set to use cuda:0
100%|██████████| 100/100 [01:07<00:00, 1.47it/s]Accuracy: 0.740
Accuracy for label Noncancer: 0.909
Accuracy for label Cancer: 0.533
```

Classification Report:

	precision	recall	f1-score	support
Noncancer	0.70	0.91	0.79	55
Cancer	0.83	0.53	0.65	45
accuracy			0.74	100
macro avg	0.77	0.72	0.72	100
weighted avg	0.76	0.74	0.73	100

Confusion Matrix:

```
[[50  5]
 [21 24]]
```

```
1 trainer.save_model(output_dir)
2 tokenizer.save_pretrained(output_dir)
```



```
('llama-3.2-fine-tuned-model/tokenizer_config.json',
'llama-3.2-fine-tuned-model/special_tokens_map.json',
'llama-3.2-fine-tuned-model/tokenizer.json')
```

```
1 # Load Model base model
2 model = AutoModelForCausalLM.from_pretrained(
3     base_model_name,
4     device_map="auto",
5     torch_dtype="float16",
6     quantization_config=bnb_config,
7 )
8
9 # Merge LoRA and base model and save
10 peft_model = PeftModel.from_pretrained(model, output_dir)
```

```

Loading checkpoint shards: 0%|          | 0/2 [00:00<?, ?it/s]
/usr/local/lib/python3.11/dist-packages/peft/tuners/lora/bnb.py:351: UserWarning: Merge lora module to 4-bit linear may
warnings.warn(
('merged-LoRA-llama-model-tokenizer/tokenizer_config.json',
 'merged-LoRA-llama-model-tokenizer/special_tokens_map.json',
 'merged-LoRA-llama-model-tokenizer/tokenizer.json')

```