

Fine-tune a small pre-trained language model (e.g., T5-small, DistilBERT) for a specific text classification or summarization task.

```
1 !pip install -q transformers datasets sentencepiece accelerate
2 import os, torch
3 os.environ["WANDB_DISABLED"]="true"
4 device=torch.device("cuda")
5
6 from datasets import load_dataset
7 from transformers import AutoTokenizer, AutoModelForSeq2SeqLM
8 from transformers import DataCollatorForSeq2Seq, TrainingArguments, Trainer
9
10 tokenizer=AutoTokenizer.from_pretrained("t5-small")
11 model=AutoModelForSeq2SeqLM.from_pretrained("t5-small").to(device)
12
13 ds=load_dataset("imdb")
14 ds=ds["train"].select(range(3000)).train_test_split(0.2)
15
16 def pre_classify(b):
17     x=["classify sentiment: "+t for t in b["text"]]
18     y=["positive" if l==1 else "negative" for l in b["label"]]
19     t=tokenizer(x, truncation=True)
20     l=tokenizer(y, truncation=True)
21     t["labels"]=l["input_ids"]
22     return t
23
24 cls=ds.map(pre_classify, batched=True, remove_columns=["text","label"])
25 collator=DataCollatorForSeq2Seq(tokenizer, model=model)
26
27 args1=TrainingArguments(
28     "t5-classification",
29     per_device_train_batch_size=4,
30     num_train_epochs=1,
31     learning_rate=3e-4,
32     logging_steps=50,
33     fp16=True
```

```
34 )
35
36 trainer1=Trainer(model=model, args=args1, train_dataset=cls["train"], eval_dataset=cls["test"], data_collator=collator)
37 trainer1.train()
38
39 def classify(x):
40     i=tokenizer("classify sentiment: "+x, return_tensors="pt").to(device)
41     o=model.generate(**i, max_new_tokens=5)
42     return tokenizer.decode(o[0], skip_special_tokens=True)
43
44 ds2=load_dataset("cnn_dailymail","3.0.0")
45 ds2=ds2["train"].select(range(2000)).train_test_split(0.1)
46
47 def pre_summarize(b):
48     x=[{"summarize": "+a for a in b['article']"}]
49     y=b["highlights"]
50     t=tokenizer(x, truncation=True)
51     l=tokenizer(y, truncation=True)
52     t["labels"]=l["input_ids"]
53     return t
54
55 sum_ds=ds2.map(pre_summarize, batched=True, remove_columns=["article","highlights"])
56
57 args2=TrainingArguments(
58     "t5-summarization",
59     per_device_train_batch_size=2,
60     num_train_epochs=1,
61     learning_rate=3e-4,
62     logging_steps=50,
63     fp16=True
64 )
65
66 trainer2=Trainer(model=model, args=args2, train_dataset=sum_ds["train"], eval_dataset=sum_ds["test"], data_collator=collator)
67 trainer2.train()
68
69 def summarize(x):
70     i=tokenizer("summarize: "+x, return_tensors="pt").to(device)
```

```
71     o=model.generate(**i, max_new_tokens=60)
72     return tokenizer.decode(o[0], skip_special_tokens=True)
73
74
75
```



```
1 print("Classification:", classify("The movie was fantastic!"))
2 print("Summarization:", summarize("The Eiffel Tower is a world-famous landmark visited by millions of tourists"))
3
```

Classification: negative

Summarization: Eiffel Tower is a world-famous landmark visited every year.

Step	Training Loss
50	2.262200
100	0.000000
150	0.000000