

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

```
1 !pip install transformers datasets peft accelerate
```

```
1 from transformers import AutoTokenizer, AutoModelForCausalLM
2 import torch
3
4 base_model = "HuggingFaceTB/SmolLM2-360M-Instruct"
5
6 tokenizer = AutoTokenizer.from_pretrained(base_model)
7 tokenizer.pad_token = tokenizer.eos_token
8
9 model = AutoModelForCausalLM.from_pretrained(
10     base_model,
11     torch_dtype=torch.float16,
12     device_map="auto"
13 )
```

```
1 import json
2 from datasets import Dataset
3
4 file_path = "/content/drive/MyDrive/Datasets/sorachio-identity-dataset.jsonl"
5
6 data = []
7 with open(file_path, "r") as f:
8     for idx, line in enumerate(f, start=1):
9         try:
10             data.append(json.loads(line))
11         except json.JSONDecodeError as e:
12             print(f"❌ JSON decode error on line {idx}: {e}")
13             print("📄 Line content:")
14             print(line)
15             break
16
17 dataset = Dataset.from_list(data)
18
19 print(dataset[:2])
20
```

```
1 from torch.nn import Linear
2
3 for name, module in model.named_modules():
4     if isinstance(module, Linear):
5         print(name)
```

```
1 from peft import LoraConfig, get_peft_model, prepare_model_for_kbit_training
2
3 lora_config = LoraConfig(
4     r=16,
5     lora_alpha=32,
6     target_modules=["q_proj", "v_proj", "k_proj", "o_proj", "gate_proj", "up_proj", "down_proj"],
7     lora_dropout=0.05,
8     bias="none",
9     task_type="CAUSAL_LM"
10 )
11
12 model = get_peft_model(model, lora_config)
13 model.print_trainable_parameters()
```

```
1 def format_prompt(example):
2     # Memisahkan berdasarkan <|im_start|> dan <|im_end|>
3     text = example["text"]
4     messages = []
5
6     # Memisahkan berdasarkan token yang ada
7     parts = text.split("<|im_end|>")
8     for part in parts:
9         if "<|im_start|>system" in part:
10             messages.append({"role": "system", "content": part.replace("<|im_start|>system\n", "").strip()})
11         elif "<|im_start|>user" in part:
```

```

12     messages.append({"role": "user", "content": part.replace("<|im_start|>user\n", "").strip()})
13     elif "<|im_start|>assistant" in part:
14         messages.append({"role": "assistant", "content": part.replace("<|im_start|>assistant\n", "").strip()})
15
16     return tokenizer.apply_chat_template(messages, tokenize=False)
17
18 def tokenize(example):
19     return tokenizer(format_prompt(example), truncation=True, padding="max_length", max_length=256)
20
21 # Tokenisasi dataset
22 tokenized_dataset = dataset.map(tokenize)

```

```

1 from transformers import TrainingArguments, Trainer, DataCollatorForLanguageModeling
2
3 training_args = TrainingArguments(
4     output_dir="./sorachio-lora",
5     per_device_train_batch_size=4,
6     gradient_accumulation_steps=2,
7     num_train_epochs=3,
8     learning_rate=2e-4,
9     warmup_ratio=0.03,
10    weight_decay=0.01,
11    fp16=True,
12    lr_scheduler_type="cosine",
13    save_strategy="epoch",
14    logging_steps=10,
15    report_to="none"
16 )
17
18 data_collator = DataCollatorForLanguageModeling(tokenizer, mlm=False)
19
20 trainer = Trainer(
21     model=model,
22     args=training_args,
23     train_dataset=tokenized_dataset,
24     data_collator=data_collator
25 )
26
27 trainer.train()

```

```

1 peft_model = model
2 model = model.merge_and_unload()
3 output_path = "/content/drive/MyDrive/Sorachio-360M-Chat/models/"
4
5 model.save_pretrained(output_path)
6 tokenizer.save_pretrained(output_path)

```