



پاییز ۱۴۰۴

درستی سنجی مدل‌های HDL

نام و نام خانوادگی دستیار آموزشی: بهزاد جنتی / مهسا ابرقانی اقدم



تمرین امتیازی

## ۱ بخش اول: مقدمه و اهداف تمرین

مقدمه: در دنیای پیچیده طراحی سخت‌افزار، اطمینان از صحت عملکرد مدل‌های توصیف سخت‌افزار (HDL) یکی از چالش‌های اصلی است. Assertion-Based Verification (ABV) یک روش قدرتمند برای این منظور است که در آن، ویژگی‌ها و رفتارهای مورد انتظار مدار به صورت خاصیت‌هایی (Properties) یا همان Assertion بیان می‌شوند. نوشتن این Assertion‌ها به صورت دستی، فرآیندی زمان‌بر و مستعد خطا است. با ظهور مدل‌های زبانی بزرگ (LLMs)، فرصت‌های جدیدی برای خودکارسازی این فرآیند به وجود آمده است. هدف تمرین:

در این تمرین، شما با استفاده از تکنیک‌های پیشرفته پردازش زبان طبیعی و یادگیری عمیق، یک مدل زبانی بزرگ (Llama 3.2 1B) را برای تولید خودکار SystemVerilog Assertion از روی کدهای HDL آموزش خواهید داد. شما در طی این فرآیند با مراحل زیر آشنا خواهید شد:

۱. راه‌اندازی محیط توسعه در Google Colab و آشنایی با پلتفرم Hugging Face.
۲. بارگذاری یک مدل زبانی بزرگ از پیش آموزش‌دیده و بهینه‌سازی آن با تکنیک کوانتیزاسیون (Quantization).
۳. بررسی و پیش‌پردازش دیتاست تخصصی VERT که شامل زوج‌های کد HDL و Assertion متناظر است.
۴. آموزش مدل با استفاده از روش نوین و بهینه LoRA (Low-Rank Adaptation) که یک تکنیک Parameter-Efficient Fine-Tuning (PEFT) است.
۵. ارزیابی عملکرد مدل قبل و بعد از فرآیند Fine-tuning و تحلیل نتایج به دست آمده.

## ۲ بخش دوم: پیش‌نیازها و آماده‌سازی محیط

قبل از شروع کدنویسی، لازم است با ابزارهای اصلی این تمرین آشنا شوید و دسترسی‌های لازم را فراهم کنید. الف) آشنایی با Hugging Face و دریافت توکن دسترسی Hugging Face یک پلتفرم پیشرو در زمینه یادگیری ماشین است که به عنوان "GitHub برای مدل‌های هوش مصنوعی" شناخته می‌شود. در این پلتفرم می‌توانید به هزاران مدل، دیتاست و ابزار دسترسی داشته باشید. برای استفاده از مدل‌های خاص (مانند مدل Llama 3.2) نیاز به

یک حساب کاربری و توکن دسترسی دارید.  
دستورالعمل:

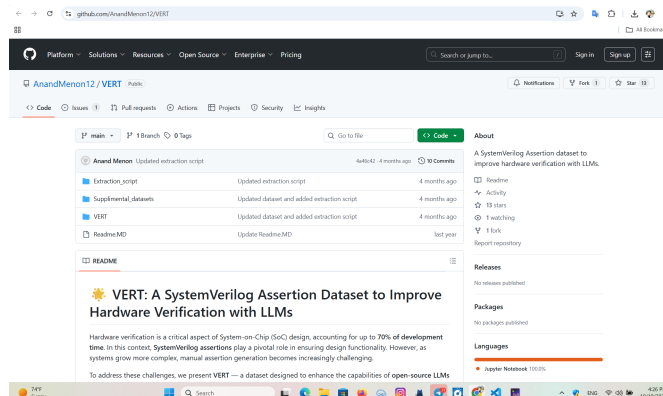


به سایت [huggingface.co](https://huggingface.co) مراجعه کرده و یک حساب کاربری ایجاد کنید. به بخش Access -> Settings Tokens در پروفایل خود بروید یا از لینک مستقیم [huggingface.co/settings/tokens](https://huggingface.co/settings/tokens) استفاده کنید. یک توکن جدید (New token) با سطح دسترسی Read ایجاد کنید. این توکن را کپی کرده و در جایی امن نگه دارید. در ادامه تمرین از آن برای احراز هویت استفاده خواهیم کرد. به صفحه مدل meta-llama/Llama-3.2-1B در Hugging Face مراجعه کرده و درخواست دسترسی (Request Access) خود را ثبت کنید. معمولاً دسترسی به سرعت تایید می‌شود.

(ب) آشنایی با دیتاست VERT

دیتاست VERT یک مجموعه داده تخصصی برای وظیفه تولید SystemVerilog Assertion از روی کدهای Verilog است. این دیتاست از طریق یک مخزن گیت‌هاب در دسترس است.  
دستورالعمل:

به مخزن گیت‌هاب دیتاست در آدرس <https://github.com/AnandMenon12/VERT> مراجعه کنید. ساختار مخزن و به خصوص فایل‌های داده را بررسی کنید.



سوال ۱: دیتاست VERT چه ساختاری دارد؟ هر نمونه داده (Sample) شامل چه فیلدهایی است و هر فیلد چه اطلاعاتی را در خود جای داده است؟ مشاهدات خود را در چند خط توضیح دهید.

**پاسخ خود را در گزارش بنویسید.**

### ۳ بخش سوم: مراحل پیاده‌سازی و کدنویسی

در این بخش، کدهای ارائه شده در نوت‌بوک را به همراه توضیحات و سوالات مرتبط دنبال کنید. شما باید تمام سلول‌های کد را اجرا کرده و به سوالات مطرح شده در سلول‌های متنی پاسخ دهید.

۱. نصب کتابخانه‌ها و راه‌اندازی محیط توضیح: در اولین قدم، کتابخانه‌های مورد نیاز پایتون را نصب می‌کنیم. این کتابخانه‌ها شامل موارد زیر هستند:

- transformers: برای بارگذاری و کار با مدل‌های از پیش آموزش دیده.
- datasets: برای مدیریت و پردازش دیتاست‌ها.
- peft: برای پیاده‌سازی تکنیک‌های Fine-tuning بهینه (مانند LoRA).
- bitsandbytes: برای کوانتیزاسیون مدل و کاهش مصرف حافظه.
- trl: برای ساده‌سازی فرآیند آموزش مدل‌های Transformer.

پس از نصب، نسخه PyTorch و وضعیت دسترسی به GPU را بررسی می‌کنیم تا از آماده بودن محیط برای اجرای محاسبات سنگین مطمئن شویم.

```
!pip install -q transformers datasets peft accelerate bitsandbytes trl sentencepiece protobuf

import torch
import os
import json
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from datasets import load_dataset, Dataset
from transformers import (
    AutoTokenizer,
    AutoModelForCausalLM,
    BitsandbytesConfig,
    TrainingArguments,
    pipeline
)
from peft import LoraConfig, get_peft_model, prepare_model_for_kbit_training
from trl import SFTTrainer
import warnings
warnings.filterwarnings('ignore')

print("\n All libraries imported successfully")
print(f"PyTorch version: {torch.__version__}")
print(f"CUDA available: {torch.cuda.is_available()}")
if torch.cuda.is_available():
    print(f"GPU: {torch.cuda.get_device_name(0)}")
    print(f"GPU memory: {torch.cuda.get_device_properties(0).total_memory / 1099.511616} GB")

66.1/66.1 MB 14.6 MB/s eta 0:00:00
564.6/564.6 kB 38.4 MB/s eta 0:00:00

✓ All libraries imported successfully
PyTorch version: 2.0.0+cu118
CUDA available: True
GPU: Tesla T4
GPU Memory: 15.83 GB
```

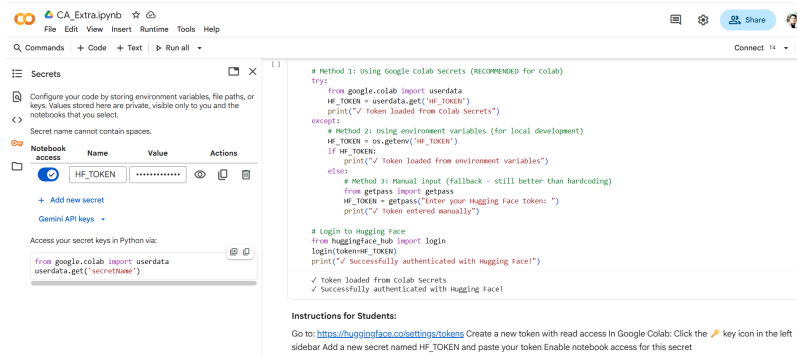
سوال ۲: چرا بررسی نسخه PyTorch و فعال بودن CUDA در ابتدای یک پروژه یادگیری عمیق اهمیت دارد؟ GPU استفاده شده در این نوت‌بوک (Tesla T4) چه خصوصیتی دارد که آن را برای این تمرین مناسب می‌کند؟

**پاسخ خود را در گزارش بنویسید**

۲. مدیریت امن توکن Hugging Face (امتیاز اضافه) توضیح: قراردادن مستقیم توکن‌های دسترسی در کد (Hardcoding) یک عمل بسیار ناامن است. در این بخش، روش صحیح برای استفاده از توکن Hugging Face در

محیط Google Colab با استفاده از Colab Secrets آموزش داده می‌شود. این روش اطلاعات حساس شما را به صورت امن ذخیره می‌کند. دستورالعمل:

در محیط Google Colab، روی آیکن کلید در نوار کناری سمت چپ کلیک کنید. یک Secret جدید با نام HF\_TOKEN ایجاد کنید. توکنی که از سایت Hugging Face کپی کردید را در قسمت Value وارد کنید. دسترسی نوت‌بوک (Notebook access) را برای این Secret فعال کنید.



سوال ۳: چرا استفاده از Colab Secrets یا متغیرهای محیطی (Environment Variables) برای ذخیره توکن، روش بهتری نسبت به نوشتن مستقیم آن در کد است؟ (در حد یک یا دو خط توضیح دهید)

پاسخ خود را در گزارش بنویسید

درخواست دسترسی به مدل Llama 3.2

مدل Llama 3.2 توسط شرکت Meta توسعه داده شده و برای استفاده از آن باید درخواست دسترسی ثبت کنید. دستورالعمل: ۱. با حساب کاربری خود وارد Hugging Face شوید. ۲. به صفحه مدل meta-llama/Llama-3.2-1B در لینک زیر بروید: <https://huggingface.co/meta-llama/Llama-3.2-1B>. ۳. شرایط استفاده را مطالعه و فرم درخواست دسترسی را تکمیل و ارسال کنید. معمولاً دسترسی در مدت کوتاهی (از چند دقیقه تا چند ساعت) از طریق ایمیل تایید می‌شود. ۴. خواسته شده: اسکرین‌شاتی از ایمیل تایید دسترسی خود را در گزارش نهایی قرار دهید.

۳. بارگذاری مدل و بهینه‌سازی با کوانتیزاسیون

توضیح: در این مرحله، مدل اصلی (meta-llama/Llama-3.2-1B) را از Hugging Face بارگذاری می‌کنیم. از آنجایی که مدل‌های زبانی بزرگ به حافظه GPU زیادی نیاز دارند، ما از یک تکنیک بهینه‌سازی به نام کوانتیزاسیون ۴ بیتی (4-bit Quantization) استفاده می‌کنیم. با استفاده از کتابخانه bitsandbytes، وزن‌های مدل که به صورت پیش‌فرض با دقت ۳۲ بیت (float32) هستند، با دقت ۴ بیت نمایش داده می‌شوند. این کار باعث کاهش چشمگیر مصرف حافظه شده و به ما اجازه می‌دهد مدل را روی یک GPU با حافظه محدود (مانند Tesla T4 با ۱۵

گیگابایت VRAM) بارگذاری و آموزش دهیم. همچنین Tokenizer مدل را بارگذاری می‌کنیم که وظیفه تبدیل متن به توکن‌های عددی (که برای مدل قابل فهم است) و برعکس را بر عهده دارد. مدل‌های زبانی بزرگ حافظه گرافیکی (VRAM) زیادی مصرف می‌کنند. مدل Llama-3.2-1B حتی در اندازه کوچک خود، در حالت عادی (32-bit precision) به بیش از ۴GB حافظه نیاز دارد. برای اینکه بتوانیم این مدل را در محیط رایگان Google Colab (که معمولاً GPU با حدود ۱۵GB VRAM دارد) آموزش دهیم، از تکنیک کوانتیزاسیون (Quantization) استفاده می‌کنیم. در اینجا، ما مدل را با دقت 4-bit بارگذاری می‌کنیم. این کار حجم مدل در حافظه را به شکل چشمگیری کاهش می‌دهد.

BitsAndBytesConfig: این کلاس به ما اجازه می‌دهد پارامترهای کوانتیزاسیون را تعریف کنیم. AutoTokenizer: وظیفه تبدیل متن به توکن (اعداد قابل فهم برای مدل) و برعکس را بر عهده دارد. AutoModelForCausalLM: کلاس اصلی برای بارگذاری مدل‌های زبانی است که برای تولید متن (Causal Language Modeling) طراحی شده‌اند.

```

1
2 # Section 2: Model and Tokenizer Loading
3 # =====
4
5 # 2.1: Define model name and quantization configuration
6 MODEL_NAME = "meta-llama/Llama-3.2-1B"
7
8 # Configure 4-bit quantization for memory efficiency
9 bnb_config = BitsAndBytesConfig(
10     load_in_4bit=True,
11     bnb_4bit_quant_type="nf4",
12     bnb_4bit_compute_dtype=torch.float16,
13     bnb_4bit_use_double_quant=True,
14 )
15
16 print(f"Loading model: {MODEL_NAME} with 4-bit quantization ...")
17
18 # 2.2: Load tokenizer

```

```

19 tokenizer = AutoTokenizer.from_pretrained(
20     MODEL_NAME,
21     trust_remote_code=True
22 )
23 # Set padding token if not present
24 if tokenizer.pad_token is None:
25     tokenizer.pad_token = tokenizer.eos_token
26     tokenizer.pad_token_id = tokenizer.eos_token_id
27
28 print("Tokenizer loaded successfully!")
29
30 # 2.3: Load model with quantization
31 model = AutoModelForCausalLM.from_pretrained(
32     MODEL_NAME,
33     quantization_config=bnb_config,
34     device_map="auto",
35     trust_remote_code=True
36 )
37
38 print("Model loaded successfully!")
39
40 # 2.4: Check GPU memory usage after loading
41 if torch.cuda.is_available():
42     print(f"\nGPU Memory allocated {torch.cuda.memory_allocated(0) / 1e9}

```

سوال ۴: در کد بالا، کلاس BitsAndBytesConfig برای کوانتیزاسیون مدل استفاده شده است. به طور خلاصه توضیح دهید هر یک از پارامترهای load\_in\_4bit, bnb\_4bit\_quant\_type و bnb\_4bit\_compute\_dtype چه نقشی دارند؟ استفاده از کوانتیزاسیون چه مزایا و معایب بالقوه‌ای دارد؟

پاسخ خود را در گزارش بنویسید

۴: بارگذاری و تحلیل دیتاست VERT

توضیح:

در این مرحله، دیتاست VERT را از مخزن گیت‌هاب آن دریافت کرده و بارگذاری می‌کنیم. پس از بارگذاری، ساختار داده‌ها را بررسی کرده و چند نمونه را برای درک بهتر وظیفه (تولید Assertion از روی کد HDL) نمایش می‌دهیم. همچنین، با رسم نمودار، آماری از ویژگی‌های دیتاست مانند طول کدها، طول Assertion‌ها و توزیع Assertion‌های همزمان (Synchronous) و غیرهمزمان (Asynchronous) به دست می‌آوریم.

```
1 # Section 3: Dataset Loading and Exploration
2 # =====
3
4 # 3.1: Clone the VERT dataset repository
5 !git clone https://github.com/AnandMenon12/VERT.git
6
7 # 3.2: Load the dataset from the JSON file
8 try:
9     with open('/content/VERT/VERT/VERT.json', 'r') as f:
10         # Each line is a JSON object
11         dataset_list = [json.loads(line) for line in f]
12         print(f"Successfully loaded {len(dataset_list)} samples from VERT")
13 except Exception as e:
14     print(f"Error loading dataset: {e}")
15     dataset_list = []
16
17 # 3.3: Split into train and test sets
18 from sklearn.model_selection import train_test_split
19 train_data, test_data = train_test_split(dataset_list, test_size=0.2, r
20
21 print(f"Training samples: {len(train_data)}")
22 print(f"Test samples: {len(test_data)}")
23
24 # 3.4: Display a few samples to understand the structure
```

```

25 print("\n" + "="*80)
26 print("SAMPLE_DATA_ANALYSIS")
27 print("="*80)
28 for i, sample in enumerate(train_data[:2], 1):
29     print(f"\n——Sample_{i}——")
30     print(f"Code:_{str(sample.get('Code', 'N/A'))[:200]}...")
31     print(f"Assertion:_{str(sample.get('Assertion', 'N/A'))[:200]}...")
32     print(f"Synchronous:_{sample.get('Synchronous', 'N/A')}")
33     print(f"Clock:_{sample.get('Clock', 'N/A')}")
34
35 # 3.5: Analyze and visualize dataset statistics
36 assertion_lengths = [len(str(s.get('Assertion', ''))).split() for s in
37 code_lengths = [len(str(s.get('Code', ''))).split() for s in train_data
38 sync_count = sum(1 for s in train_data if s.get('Synchronous'))
39
40 plt.figure(figsize=(15, 5))
41 plt.subplot(1, 2, 1)
42 sns.histplot(assertion_lengths, bins=30, kde=True, color='blue')
43 plt.title('Distribution_of_Assertion_Lengths_(in_words)')
44 plt.subplot(1, 2, 2)
45 sns.histplot(code_lengths, bins=30, kde=True, color='green')
46 plt.title('Distribution_of_Code_Lengths_(in_words)')
47 plt.show()
48
49 print(f"\nPercentage_of_Synchronous_Assertions_in_training_data:_{sync_

```

سوال ۵: الف) بر اساس نمودارهای توزیع طول کد و Assertion، چه نتیجه‌ای در مورد ماهیت دیتاست VERT می‌توان گرفت؟ آیا طول Assertion‌ها با طول کدها رابطه مستقیمی دارد؟ ب) چرا تقسیم داده‌ها به دو مجموعه آموزش (Train) و آزمون (Test) یک مرحله ضروری در پروژه‌های یادگیری ماشین است؟

پاسخ خود را در گزارش بنویسید

۵: فرمت‌بندی داده‌ها برای آموزش

توضیح:

برای اینکه مدل Llama 3.2 یاد بگیرد که چگونه از یک کد HDL، یک Assertion تولید کند، باید داده‌ها را در یک فرمت خاص به نام InstructionFollowing آماده کنیم. در این فرمت، ما یک “دستورالعمل” (شامل توضیح وظیفه و کد ورودی) و یک “پاسخ” (Assertion مطلوب) ایجاد می‌کنیم. این کار به مدل کمک می‌کند تا بفهمد ورودی و خروجی مورد انتظار چه ساختاری دارند. تابع `format_instruction` این وظیفه را انجام می‌دهد.

```

1 # Section 4: Data Formatting for Instruction Tuning
2 # =====
3
4 def format_instruction(sample):
5     """Formats a sample into the LLaMA 3.2 instruction format."""
6     code = sample.get('Code', '')
7     assertion = sample.get('Assertion', '')
8     context = f"This is a synchronous assertion that executes at {sample.get('clock', 'clock')}
9
10    # LLaMA 3.2 instruction format
11    return f"""<|begin_of_text|><|start_header_id|>system<|end_header_id|>
12
13    You are an expert in hardware verification. Generate SystemVerilog assertions for the following
14
15    ### HDL Code:
16    {code}
17
18    ### Assertion:
19    {assertion}<|eot_id|>"""
20
21 # Format the datasets
22 train_texts = [format_instruction(sample) for sample in train_data]
23 test_texts = [format_instruction(sample) for sample in test_data]
```

```

24
25 train_dataset = Dataset.from_dict({"text": train_texts})
26 test_dataset = Dataset.from_dict({"text": test_texts})
27
28 print(" _Datasets_formatted_for_instruction_tuning.")
29 print("\n——Example_of_a_Formatted_Training_Sample——")
30 print(train_texts[0])

```

سوال ۶: در تابع `format_instruction`، از توکن‌های خاصی مانند `<|begin_of_text|>`، `<|start_header_id|>` و `<|eot_id|>` استفاده شده است. نقش این توکن‌های خاص در آموزش مدل‌های زبانی چیست و چرا استفاده از فرمت صحیح برای یک مدل خاص (در اینجا Llama 3.2) اهمیت دارد؟

پاسخ خود را در گزارش بنویسید

## ۴ بخش چهارم: آموزش مدل با استفاده از LoRA

۶. پیکربندی LoRA برای FineTuning بهینه

توضیح: Full Fine-Tuning یا آموزش کامل یک مدل زبانی بزرگ، به ده‌ها یا صدها گیگابایت حافظه GPU نیاز دارد. برای حل این مشکل، از روش‌های (Parameter-Efficient Fine-Tuning (PEFT استفاده می‌کنیم. یکی از محبوب‌ترین این روش‌ها، LoRA (LowRank Adaptation) است.

در تکنیک LoRA، به جای به‌روزرسانی تمام پارامترهای مدل اصلی (که میلیاردها پارامتر هستند)، وزن‌های اصلی مدل را فریز (Freeze) کرده و فقط ماتریس‌های کوچکی با رتبه پایین (LowRank) را که به لایه‌های مشخصی از مدل (معمولاً لایه‌های Attention) اضافه می‌شوند، آموزش می‌دهیم. این کار باعث می‌شود تعداد پارامترهای قابل آموزش به شدت کاهش یابد (معمولاً کمتر از 1٪ کل پارامترها) و فرآیند Finetuning با منابع سخت‌افزاری محدود نیز ممکن شود. در کد زیر، ما LoRA را برای لایه‌های کلیدی مدل (`q_proj`، `k_proj`، `v_proj`، و غیره) پیکربندی می‌کنیم.

```

1 # Prepare model for k-bit training
2 model = prepare_model_for_kbit_training(model)
3
4 # Configure LoRA for parameter-efficient fine-tuning
5 lora_config = LoraConfig(

```

```

6         r=16,                                # Rank of the update matrices
7         lora_alpha=32,                        # Scaling factor for LoRA weights
8         target_modules=[
9             "q_proj", "k_proj", "v_proj", "o_proj",
10            "gate_proj", "up_proj", "down_proj"
11        ],
12        lora_dropout=0.05,                    # Dropout probability for LoRA layers
13        bias="none",
14        task_type="CAUSAL_LM"
15    )
16
17    # Apply LoRA to the model
18    model = get_peft_model(model, lora_config)
19
20    print("\u2713 LoRA configured successfully!")
21
22    # Print the percentage of trainable parameters
23    model.print_trainable_parameters()
24
25    # Calculate parameter efficiency
26    trainable_params = sum(p.numel() for p in model.parameters() if p.requires_grad)
27    total_params = sum(p.numel() for p in model.parameters())
28    trainable_percentage = 100 * trainable_params / total_params
29    print(f"\n\u25b6 Percentage of trainable parameters: {trainable_percentage}%")

```

سوال ۷: الف) تکنیک LoRA چیست و چگونه به کاهش چشمگیر مصرف حافظه و منابع محاسباتی در فرآیند Finetuning کمک می‌کند؟ ب) خروجی دستور `model.print_trainable_parameters()` چه چیزی را نشان می‌دهد؟ چرا تعداد پارامترهای قابل آموزش (trainable params) بسیار کمتر از کل پارامترهای مدل است؟

پاسخ خود را در گزارش بنویسید

۷. پیکربندی و اجرای فرآیند آموزش

توضیح: در این مرحله، هایپرپارامترهای فرآیند آموزش را با استفاده از کلاس `TrainingArguments` مشخص می‌کنیم. این پارامترها شامل تعداد `epoch`، اندازه `batch`، نرخ یادگیری (`learning rate`) و  $\beta$  هستند. یکی از پارامترهای مهم در اینجا `gradient_accumulation_steps` است. این تکنیک به ما اجازه می‌دهد تا یک `batch` size بزرگتر را شبیه‌سازی کنیم. به این صورت که گرادینت‌ها در چندین مرحله (در اینجا ۴ مرحله) محاسبه و تجمیع می‌شوند و تنها پس از آن، وزن‌های مدل یک بار به‌روزرسانی می‌شوند. این کار برای آموزش روی GPUهای با حافظه محدود بسیار مفید است. در نهایت، با استفاده از کلاس `SFTTrainer` (Supervised Fine-tuning Trainer) از کتابخانه `trl`، فرآیند آموزش را آغاز می‌کنیم. این کلاس به طور خاص برای `Instruction Fine-Tuning` مدل‌های `Transformer` بهینه‌سازی شده است.

```

1 # Training Arguments Configuration
2 training_args = TrainingArguments(
3     output_dir=". / llama-vert-lora-checkpoints",
4     num_train_epochs=1,                                # We will train for one full p
5     per_device_train_batch_size=2,                       # Batch size per GPU
6     gradient_accumulation_steps=4,                       # Accumulate gradients over 4
7     gradient_checkpointing=True,                         # Use gradient checkpointing to
8     learning_rate=2e-4,                                  # Learning rate
9     fp16=True,                                           # Use 16-bit precision for tra
10    max_grad_norm=0.3,                                   # Gradient clipping
11    warmup_ratio=0.03,                                   # Warmup steps ratio
12    lr_scheduler_type="cosine",                          # Cosine learning rate schedul
13    logging_steps=25,                                    # Log training progress every
14    save_strategy="epoch",                               # Save model at the end of each
15    report_to="none"
16 )
17
18 # Initialize the SFTTrainer
19 trainer = SFTTrainer(
20     model=model,
21     train_dataset=train_dataset,

```

```

22     peft_config=lora_config ,
23     dataset_text_field="text" ,
24     max_seq_length=512,                # Truncate long sequences
25     tokenizer=tokenizer ,
26     args=training_args ,
27 )
28
29 print("\u2713 Trainer initialized successfully!")
30
31 # Start the training process
32 print("\n" + "="*80)
33 print("\ud83d\ude80 STARTING FINE-TUNING PROCESS...")
34 print("="*80)
35 trainer.train()
36 print("\u2713 Training completed successfully!")
37
38 # Save the final fine-tuned model
39 final_model_path = "./llama-vert-lora-final"
40 trainer.save_model(final_model_path)
41 print(f"\u2713 Final model saved to: {final_model_path}")

```

سوال ۸: الف) در لاگ‌های خروجی فرآیند آموزش، مقدار Loss چه چیزی را نشان می‌دهد و چرا انتظار داریم این مقدار در طول زمان کاهش یابد؟

ب) با توجه به مقادیر `per_device_train_batch_size` و `gradient_accumulation_steps`، اندازه batch مؤثر (Effective Batch Size) در این آموزش چقدر است؟ این تکنیک چه مزیتی دارد؟

پاسخ خود را در گزارش بنویسید:

## ۵ بخش پنجم: ارزیابی و تحلیل نتایج

حالا زمان آن رسیده که ببینیم فرآیند Fine-tuning چقدر موفقیت‌آمیز بوده است. در این بخش، عملکرد مدل را پس از آموزش ارزیابی کرده و آن را با عملکرد اولیه مدل مقایسه می‌کنیم.

۸. ارزیابی مدل Fine-tune شده توضیح: در این بخش، ابتدا مدل آموزش‌دیده را به همراه آداپتورهای LoRA بارگذاری می‌کنیم. سپس یک pipeline برای تولید متن ایجاد کرده و از آن برای تولید Assertion روی همان ۱۰ نمونه از مجموعه آزمون که در بخش ارزیابی اولیه استفاده کردیم، بهره می‌گیریم. این کار به ما اجازه می‌دهد تا مقایسه‌ای مستقیم بین خروجی‌های مدل قبل و بعد از آموزش داشته باشیم.

```

1 # Import necessary libraries for evaluation
2 from peft import PeftModel
3 from transformers import AutoModelForCausalLM, AutoTokenizer, pipeline
4 import torch
5
6 # Path to the base model and the saved LoRA adapter
7 base_model_name = "meta-llama/Llama-3.2-1B"
8 adapter_path = "./llama-vert-lora-final"
9
10 print("\n\ud83d\udd0d Loading base model and merging LoRA adapter for \ud83d\udd0d")
11
12 # Load the base model in 4-bit
13 base_model = AutoModelForCausalLM.from_pretrained(
14     base_model_name,
15     quantization_config=bnb_config,
16     device_map="auto",
17     trust_remote_code=True
18 )
19
20 # Load the LoRA model and merge it with the base model
21 model_finetuned = PeftModel.from_pretrained(base_model, adapter_path)
22 model_finetuned = model_finetuned.merge_and_unload() # Merge weights for

```

```
23
24 tokenizer = AutoTokenizer.from_pretrained(base_model_name, trust_remote_code=True)
25 if tokenizer.pad_token is None:
26     tokenizer.pad_token = tokenizer.eos_token
27
28 print("\u2713 Fine-tuned model loaded and merged successfully!")
29
30 # Create a new pipeline for the fine-tuned model
31 pipe_finetuned = pipeline(
32     "text-generation",
33     model=model_finetuned,
34     tokenizer=tokenizer,
35     max_new_tokens=256,
36     do_sample=True,
37     temperature=0.7,
38     top_p=0.95
39 )
40
41 print("\n\u2192 Generating assertions with the FINE-TUNED model:")
42
43 # Run evaluation on the same test samples
44 # (The 'generate_assertion' function is defined in the pre-fine-tuning script)
45 finetuned_results = []
46 for i, sample in enumerate(test_data[:10], 1):
47     print(f"\n{'='*80}")
48     print(f"Test Sample {i}")
49     print(f"{'='*80}")
50     generated = generate_assertion(sample, pipe_finetuned)
51     print(f"\u2192 Ground Truth: \n\n{sample['Assertion'][:200]}")
```

```
52 print ( f "\n\ud83d\udd2e_Model_Generated_(AFTER_Fine-tuning):\n\n{ g
53 finetuned_results.append({ 'ground_truth': sample[ 'Assertion' ], 'ge
```

سوال ۹: نتایج تولید شده توسط مدل Fine-tune شده را با نتایج مدل پایه (قبل از آموزش) مقایسه کنید. حداقل دو نمونه را انتخاب کرده و به صورت کیفی توضیح دهید که عملکرد مدل در چه جنبه‌هایی (مانند رعایت سینتکس SystemVerilog، ارتباط منطقی با کد HDL، تولید ساختارهای صحیح property و endproperty) بهبود یافته است.

بخش ششم: تحلیل نهایی و جمع‌بندی در این بخش پایانی، شما باید بر اساس تمام مراحل انجام شده، یک تحلیل کلی از پروژه ارائه دهید.

سوال ۱۱ (تحلیل کلی): کل فرآیند انجام شده در این تمرین را در چند پاراگراف خلاصه کنید. هدف چه بود، چه مراحل طی شد و نتیجه اصلی چه بود؟

**پاسخ خود را در گزارش بنویسید:**

سوال ۱۲ (نقاط ضعف و محدودیت‌ها): با بررسی خروجی‌های مدل Fine-tune شده، فکر می‌کنید بزرگترین نقاط ضعف آن چیست؟ آیا مدل توانایی تولید Assertion‌های پیچیده (مانند Assertion‌های مبتنی بر دنباله یا Sequence-based) را دارد؟ در چه مواردی بیشتر دچار خطا می‌شود؟

**پاسخ خود را در گزارش بنویسید:**

سوال ۱۳ (پیشنهادات برای بهبود): بر اساس مشاهدات خود و مفاهیمی که در کلاس آموخته‌اید، حداقل سه راهکار مشخص برای بهبود عملکرد این مدل پیشنهاد دهید. این راهکارها می‌تواند شامل استفاده از مدل‌های متفاوت، تغییر در داده‌ها یا روش پیش‌پردازش، تنظیم هایپرپارامترها، یا استفاده از معیارهای ارزیابی پیشرفته‌تر باشد. برای هر پیشنهاد، دلیل خود را نیز به طور خلاصه بیان کنید.

**پاسخ خود را در گزارش بنویسید:**