

# Programming Assignment 1 - LLM Fine Tuning

## Assignment Overview

### **Goal: Fine-Tuning LLaMA on 1 GPU with Memory Optimizations**

#### **Dataset**

Climate documents dataset of IPCC reports and climate change and AI publications from the last 5 years.

- PDF file format
- On the cloud burst compute file system: /scratch/BDML25SP/

Students can use random 90% files from the dataset for fine-tuning. Note that you need to extract the text from the data and create txt files as part of preprocessing.

#### **Pretrained Model**

- LLaMA 3B model
- On the cloud burst compute file system: /scratch/BDML25SP/

#### **Key Focus**

We will be focusing on memory optimizations, for example

- Low Rank Adaptation (LoRA)
- Mixed Precision Fine-Tuning, Quantization
- Gradient Accumulation and Checkpointing

The goal of the assignment is to increase the batch size of the training process as much as possible.

#### **Deliverables**

1. A report documenting:
  - a. Memory optimization strategies used.
  - b. Training performance (maximum batch size) and evaluation results.
  - c. Step by step guide on how to run the training code.
2. Code access on HPC

#### **Evaluation**

Compute the perplexity metric on the remaining 10% of the dataset. The assignment will be evaluated primarily on the basis of how memory and time efficient the fine-tuning code is, and the final perplexity score will not hold much weight.

## Data Processing

**Digital PDFs** have **embedded text**, making them directly readable by computers. Text can be extracted easily using tools like **PyPDF2** ([docs](#)), **pdfplumber** ([docs](#)), or **PyMuPDF** ([docs](#)) in Python.

Example code

```
1  from PyPDF2 import PdfReader
2
3  # Read the PDF file
4  pdf_path = '/path/to/data/file.pdf'
5  reader = PdfReader(pdf_path)
6
7  # Extract text from the first page
8  first_page = reader.pages[0].extract_text()
9  first_page[:1000] # Display the first 1000 characters as a preview
```

To split the dataset, you can create a list of filenames in the folder and split the list in 9:1 for train and test sets.

## Low Rank Adaptation (LoRA)

**LoRA (Low-Rank Adaptation)** is a **parameter-efficient fine-tuning (PEFT) technique** designed to fine-tune large language models (LLMs) like LLaMA, GPT, or BERT with **significantly fewer trainable parameters and lower memory requirements**. It adds **trainable low-rank matrices** to **pre-trained model weights** instead of updating the entire model, making it faster and more memory-efficient.

You can fine-tune LLaMA-3B with LoRA using Hugging Face's **peft** library.

Example code

```

1  from transformers import AutoModelForCausalLM, AutoTokenizer
2  from peft import get_peft_model, LoraConfig, TaskType
3
4  # Load pre-trained LLaMA model
5  model_name = "path/to/model"
6  model = AutoModelForCausalLM.from_pretrained(model_name)
7  tokenizer = AutoTokenizer.from_pretrained(model_name)
8
9  # Define LoRA config
10 lora_config = LoraConfig(
11     task_type=TaskType.CAUSAL_LM, # Language modeling
12     r=8, # Low-rank matrix dimension
13     lora_alpha=32, # Scaling factor
14     lora_dropout=0.1, # Dropout rate
15 )
16
17 # Apply LoRA
18 lora_model = get_peft_model(model, lora_config)
19 lora_model.print_trainable_parameters()
20

```

## Precision Optimization

Using **FP16 (Half-Precision)** significantly reduces memory usage while maintaining training stability.

Example code

```

1  from transformers import TrainingArguments, Trainer
2
3  training_args = TrainingArguments(
4      output_dir="./llama-finetune",
5      ...
6      fp16=True, # Enables FP16 (Half-Precision)
7      bf16=False, # Use BF16 (Better for A100 GPUs)
8  )
9
10
11 trainer = Trainer(
12     model=model, # Fine-tuned LLaMA model
13     args=training_args, # Training arguments
14     train_dataset=train_dataset, # Training data
15     ...
16     tokenizer=tokenizer, # Tokenizer
17 )
18
19
20 trainer.train() # Start training
21

```

QLoRA allows efficient fine-tuning of LLMs **without full model updates** by using **4-bit quantization + LoRA (Low-Rank Adaptation)**.

## Dependencies:

`pip install bitsandbytes transformers accelerate peft`

## Example code

```
1  from transformers import AutoModelForCausalLM, AutoTokenizer, BitsAndBytesConfig
2  from peft import LoraConfig, get_peft_model
3
4  # Load LLaMA-3B with 4-bit quantization
5  model_name = "path/to/model"
6  quantization_config = BitsAndBytesConfig(
7      load_in_4bit=True, # Enables 4-bit quantization
8      bnb_4bit_compute_dtype=torch.float16, # Compute in FP16
9      bnb_4bit_use_double_quant=True, # Further memory optimization
10 )
11
12 # Load model with quantization
13 model = AutoModelForCausalLM.from_pretrained(model_name, quantization_config=quantization_config)
14 tokenizer = AutoTokenizer.from_pretrained(model_name)
15
16 # Apply LoRA on quantized model
17 lora_config = LoraConfig(
18     r=8, # Low-rank dimension
19     lora_alpha=32,
20     target_modules=["q_proj", "v_proj"], # Apply LoRA to attention layers
21     lora_dropout=0.1,
22     bias="none"
23 )
24 lora_model = get_peft_model(model, lora_config)
25 lora_model.print_trainable_parameters()
```

# Gradient Accumulation and Checkpointing

**Gradient Accumulation** – Simulate larger batch sizes without exceeding memory.

**Gradient Checkpointing** – Save memory by recomputing activations instead of storing them.

## Example code

```
1  training_args = TrainingArguments(
2      ...
3      per_device_train_batch_size=2, # Small batch size per GPU
4      gradient_accumulation_steps=8, # Simulate batch size of 16
5      ...
6  )
```

## Example code

```
1  from transformers import AutoModelForCausalLM
2
3  model = AutoModelForCausalLM.from_pretrained("model/path")
4  model.gradient_checkpointing_enable() # Enable memory optimization
5
```

# Evaluation

On the test set, you can evaluate the fine-tuned model by using the perplexity metric.

**Perplexity (PPL)** is a measure of how well a language model predicts a given dataset on the next token prediction task. It is commonly used to evaluate fine-tuned language models.

Mathematically, perplexity is the exponentiated average negative log-likelihood of the model predictions:

$$PPL = e^{\left(-\frac{1}{N} \sum_{i=1}^N \log P(w_i)\right)}$$

Where,  $P(w_i)$  is the probability assigned by the model to the  $i$ th word, and  $N$  is the number of words in the dataset.

**Lower perplexity means better predictions** (more confident and accurate). Higher perplexity indicates poor model performance (more uncertain predictions).

You can evaluate perplexity on your fine-tuned LLaMA model using the Hugging Face **transformers** library.

Example code

```
1  import torch
2  import math
3  from transformers import AutoModelForCausalLM, AutoTokenizer
4
5  # Load fine-tuned LLaMA model (Replace with your model path)
6  model_name = "your-finetuned-llama"
7  tokenizer = AutoTokenizer.from_pretrained(model_name)
8  model = AutoModelForCausalLM.from_pretrained(model_name)
9
10 # Sample evaluation text (use held-out dataset)
11 evaluation_text = ""
12 Climate change is caused by an increase in greenhouse gases such as CO2.
13 ""
14 tokens = tokenizer(evaluation_text, return_tensors="pt", truncation=True, padding=True)
15
16 # Compute loss
17 with torch.no_grad():
18     outputs = model(**tokens, labels=tokens["input_ids"])
19     loss = outputs.loss.item()
20
21 # Compute perplexity
22 perplexity = math.exp(loss)
23 print(f"Perplexity: {perplexity:.2f}")
24
```

Interpreting perplexity score:

Perplexity Score	Interpretation
$\leq 10$	Very good model (accurate predictions)
10 - 50	Decent model (still useful, but can be improved)
> 100	Poor performance (struggles with predictions)
> 1000	Model is guessing almost randomly