

ADVANCED FINE-TUNING STRATEGIES FOR MATHEMATICAL REASONING

Filbert Hamijoyo

The Chinese University of Hong Kong, Shenzhen
Shenzhen, China
filberthamijoyo@icloud.com

ABSTRACT

This study investigates the impact of advanced fine-tuning strategies on language models for mathematical reasoning within resource-constrained environments. We compare two contrasting models—Mistral-7B and TinyLlama-1.1B—to understand how model size, data quality, and training techniques affect performance outcomes. Using multi-source mathematical datasets and Parameter-Efficient Fine-Tuning (PEFT) with optimized Low-Rank Adaptation (LoRA), we demonstrate the critical role of both model selection and training methodology. Our results show a significant baseline capability gap, with Mistral-7B achieving 46% baseline accuracy compared to TinyLlama’s 0%, highlighting the advantage of larger models. However, through our enhanced fine-tuning pipeline, the TinyLlama model achieved a remarkable improvement to 40% accuracy, demonstrating that even small models can develop specialized capabilities with proper training techniques. The Mistral-7B model improved to 54% accuracy after fine-tuning, showing that both large and small models benefit from our enhanced methodology. These findings provide valuable insights into optimizing fine-tuning strategies when working with limited computational resources, emphasizing the importance of data diversity, prompt engineering, and parameter-efficient adaptation techniques.

1 INTRODUCTION

Research Topic This research investigates advanced fine-tuning strategies for mathematical reasoning in resource-constrained environments, focusing on optimizing training approaches for both large and small language models. Specifically, we examine two contrasting models: Mistral-7B-SFT-Beta (7 billion parameters) and TinyLlama-1.1B-Chat (1.1 billion parameters). Our central research question explores how data quality, prompt engineering, and parameter-efficient adaptation methods can overcome computational limitations to achieve meaningful improvements in specialized reasoning tasks.

Task Importance Fine-tuning large language models for specialized tasks like mathematical reasoning is increasingly important as organizations seek to develop domain-specific AI applications. However, the computational resources required for effective fine-tuning present a significant barrier to entry for many researchers and smaller organizations. This challenge raises important questions about the accessibility of advanced AI capabilities and whether innovative training approaches can bridge the gap between model size and performance.

Mathematical reasoning represents a particularly valuable benchmark for several reasons:

- It requires structured, step-by-step logical thinking
- Performance is objectively measurable through answer accuracy
- It has practical applications across educational, scientific, and business domains
- It tests a model’s ability to follow instructions and format outputs appropriately

Understanding how to optimize fine-tuning strategies for mathematical reasoning provides valuable insights that can generalize to other complex reasoning tasks, potentially democratizing access to advanced AI capabilities.

Approach and Achievements Our approach involved parallel implementation of fine-tuning pipelines for two different-sized models, with particular emphasis on developing advanced techniques for memory efficiency and performance optimization:

- We implemented a progressive fallback strategy for model loading that enables working with larger models (7B parameters) even on consumer-grade GPUs with limited VRAM
- We developed enhanced data processing techniques that combine multiple mathematical datasets for more robust training signals
- We optimized the Low-Rank Adaptation (LoRA) configuration with higher rank and alpha values to improve the expressive capacity of parameter-efficient fine-tuning
- We implemented sophisticated prompt engineering to guide models toward structured reasoning

Both implementations utilized the Parameter-Efficient Fine-Tuning (PEFT) approach with Low-Rank Adaptation (LoRA) to minimize memory requirements while still allowing meaningful parameter updates.

Our key findings reveal:

- Base model performance varied dramatically by model size, with Mistral-7B achieving 46% accuracy compared to TinyLlama-1.1B's 0% on the GSM8K benchmark
- The enhanced fine-tuning pipeline produced remarkable improvements in the TinyLlama model, boosting its accuracy to 40%, a gain of 40 percentage points
- The Mistral-7B model improved to 54% accuracy after fine-tuning, showing that our methods benefit larger models as well
- Advanced memory optimization techniques enabled successful training and evaluation of larger models within the constraints of consumer-grade GPUs
- Multi-source training data and improved prompt engineering proved particularly valuable for developing robust mathematical reasoning capabilities

These findings highlight both the challenges and opportunities in fine-tuning LLMs under resource constraints, demonstrating that methodological innovations can significantly narrow the performance gap between small and large models for specialized tasks.

2 EXPERIMENT DESIGN

Task Definition Our experiment focuses on enhancing mathematical reasoning capabilities in language models through fine-tuning. Specifically, we define mathematical reasoning in this context as the ability to:

- Parse natural language descriptions of mathematical problems
- Identify the key variables and relationships in the problem
- Apply appropriate mathematical operations and formulas
- Execute calculations accurately through a chain-of-thought process
- Present solutions in a structured, step-by-step manner
- Arrive at the correct numerical answer

Success is measured primarily through answer accuracy on a held-out evaluation set from the GSM8K benchmark, which contains grade school math word problems requiring multi-step reasoning. We specifically instructed models to provide step-by-step solutions and clearly indicate final answers after a "" marker to facilitate consistent evaluation.

Experiment Design We designed parallel experiments for two different-sized models to investigate the impact of model size and computational resources on fine-tuning outcomes. Our overall approach included:

1. Model Selection:

- **Mistral-7B-SFT-Beta:** A 7 billion parameter model known for strong reasoning capabilities
- **TinyLlama-1.1B-Chat:** A much smaller 1.1 billion parameter model designed for efficiency

2. Dataset Preparation:

- Combined multiple datasets including PrimeIntellect/NuminaMath-QwQ-CoT-5M, GSM8K, and reasoning-machines/gsm-hard to create a diverse training corpus
- Implemented improved data processing techniques including chunking for memory efficiency
- Enhanced prompting templates to encourage step-by-step reasoning

3. Training Methodology:

- Implemented Parameter-Efficient Fine-Tuning (PEFT) with Low-Rank Adaptation (LoRA)
- Optimized LoRA configuration with increased rank and alpha values
- Utilized advanced memory optimization techniques including gradient checkpointing, 4-bit quantization, and CPU offloading
- Employed progressive fallback strategies for model loading to handle memory constraints

4. Evaluation Protocol:

- Evaluated models on a consistent set of 50 examples from the GSM8K benchmark
- Implemented robust answer extraction techniques to handle various formats
- Measured performance using exact match accuracy after normalization
- Conducted detailed analysis of model improvements and failure cases

By implementing parallel experiments with different-sized models, we aimed to isolate the impact of model size and computational requirements while exploring how advanced training techniques could bridge performance gaps.

3 CODE IMPLEMENTATION

Our implementation consisted of two separate Python scripts, each targeting a different model architecture but following similar methodological approaches. Both implementations were developed and executed in Google Colab with a T4 GPU, which presented significant memory constraints that necessitated careful optimization.

3.1 COMMON IMPLEMENTATION COMPONENTS

Both implementations shared several key components:

```
1 # Environment setup
2 !pip install -q h5py typing-extensions wheel
3 !pip install -q -U bitsandbytes
4 !pip install -q -U git+https://github.com/huggingface/transformers
  .git
5 !pip install -q -U git+https://github.com/huggingface/peft.git
6 !pip install -q -U git+https://github.com/huggingface/accelerate.
  git
```

```

7  !pip install -q datasets
8
9  # Memory management function
10 def clear_memory():
11     """Improved memory clearing function"""
12     gc.collect()
13     torch.cuda.empty_cache()
14     if torch.cuda.is_available():
15         torch.cuda.synchronize()

1  # Dataset loading functions
2  def load_combined_math_datasets(max_examples=10000):
3     """IMPROVED: Load multiple math datasets for better training signal"""
4     print(f"Loading multiple math datasets (up to {max_examples}
5           examples total)...")
6     datasets = []
7
8     # 1. Load NuminaMath dataset (primary source)
9     try:
10         numina_dataset = load_dataset("PrimeIntellect/NuminaMath-
11                                       QwQ-CoT-5M")
12         # Sample more examples
13         sampled_numina = numina_dataset["train"].shuffle(seed=42).
14             select(
15                 range(min(5000, len(numina_dataset["train"])))
16             )
17         datasets.append(sampled_numina)
18         print(f"Added {len(sampled_numina)} examples from
19               NuminaMath")
20     except Exception as e:
21         print(f"Error loading NuminaMath: {e}")
22
23     # 2. Add GSM8K training data for better alignment with
24     #    evaluation
25     try:
26         gsm8k_train = load_dataset("gsm8k", "main")["train"]
27         # Prioritize GSM8K examples by taking more of them
28         sampled_gsm8k = gsm8k_train.shuffle(seed=42).select(
29             range(min(3000, len(gsm8k_train)))
30         )
31         datasets.append(sampled_gsm8k)
32         print(f"Added {len(sampled_gsm8k)} examples from GSM8K
33               train set")
34     except Exception as e:
35         print(f"Error loading GSM8K: {e}")
36
37     # 3. Add examples from reasoning-machines dataset for
38     #    diversity
39     try:
40         cot_dataset = load_dataset("reasoning-machines/gsm-hard",
41                                    split="train")
42         sampled_cot = cot_dataset.shuffle(seed=42).select(range(
43             min(1000, len(cot_dataset))))
44         datasets.append(sampled_cot)
45         print(f"Added {len(sampled_cot)} examples from reasoning-
46               machines/gsm-hard")
47     except Exception as e:
48         print(f"Error loading reasoning-machines dataset: {e}")

```

```

39
40     # Combine all datasets and limit total size
41     combined_data = concatenate_datasets(datasets)
42     combined_data = combined_data.shuffle(seed=42)
43     if len(combined_data) > max_examples:
44         combined_data = combined_data.select(range(max_examples))
45
46     print(f"Final combined dataset size: {len(combined_data)}
47         examples")
48     return combined_data

```

```

1  # Enhanced answer extraction function
2  def extract_answer(response, question_type="math-reasoning"):
3      """Extract the final answer with improved pattern matching"""
4      if question_type == "math-reasoning":
5          # Try multiple extraction strategies in order of
6              reliability
7
8          # 1. Find answer after #### marker (most reliable)
9          hash_match = re.search(r"#\{3,\} s*([-+]?\d*\.\d+)",
10                                response, re.DOTALL)
11          if hash_match:
12              return hash_match.group(1).strip()
13
14          # 2. Look for explicit "The answer is X" pattern
15          answer_match = re.search(r"(?: the\s+answer\s+is | final\s+
16                                  answer\s+is | answer\s+[=:]) s*([-+]?\d*\.\d+)",
17                                  response.lower(), re.DOTALL)
18          if answer_match:
19              return answer_match.group(1).strip()
20
21          # 3. Look for "Therefore" pattern
22          therefore_match = re.search(r"(?: therefore | thus | hence | so)
23                                     ,? s*([-+]?\d*\.\d+)",
24                                     response.lower(), re.DOTALL)
25          if therefore_match:
26              return therefore_match.group(1).strip()
27
28          # 4. Fallback to last number in the response
29          numbers = re.findall(r"([-+]?\d*\.\d+)", response)
30          if numbers:
31              return numbers[-1]
32
33          return response.strip()
34     else:
35         # For other question types, just return the last sentence
36             or phrase
37
38         response = response.strip()
39         sentences = response.split(".")
40         if sentences:
41             last_sentence = sentences[-1].strip()
42             if len(last_sentence) > 50:
43                 last_sentence = last_sentence[-50:].strip()
44             return last_sentence
45
46         return response.strip()

```

3.2 MISTRAL-7B IMPLEMENTATION

The implementation for the Mistral-7B model was enhanced with advanced memory optimization techniques to operate within the constraints of a Google Colab T4 GPU:

```

1 def load_model_and_tokenizer(model_name, load_in_4bit=True):
2     """Load model and tokenizer with optimized memory settings"""
3     print(f"Loading {model_name} with 4-bit quantization...")
4
5     # Configure quantization with additional CPU offloading
6     # parameters
7     quantization_config = BitsAndBytesConfig(
8         load_in_4bit=load_in_4bit,
9         bnb_4bit_compute_dtype=torch.float16,
10        bnb_4bit_use_double_quant=True,
11        bnb_4bit_quant_type="nf4"
12    )
13
14    tokenizer = AutoTokenizer.from_pretrained(model_name)
15
16    # Handle tokenizer peculiarities
17    if not tokenizer.pad_token:
18        tokenizer.pad_token = tokenizer.eos_token
19
20    print("Loading with advanced memory optimization...")
21
22    # Check available GPU memory and determine if we need disk
23    # offloading
24    try:
25        free_in_GB = torch.cuda.get_device_properties(0).
26            total_memory / 1e9
27        max_memory = {0: f"{int(free_in_GB * 0.85)}GB"}
28        print(f"GPU memory available: {free_in_GB:.2f} GB,
29            allocating: {max_memory}")
30    except:
31        max_memory = None
32        print("Could not determine GPU memory, using default
33            allocation")
34
35    # IMPROVED: Progressive fallback strategy for model loading
36    try:
37        # First try with 4-bit quantization
38        model = AutoModelForCausalLM.from_pretrained(
39            model_name,
40            device_map="auto",
41            quantization_config=quantization_config,
42            torch_dtype=torch.float16,
43            offload_folder="offload_folder",
44            max_memory=max_memory,
45            offload_state_dict=True # Offload weights to CPU when
46                not in use
47        )
48    except Exception as e:
49        print(f"Initial loading attempt failed: {e}")
50        print("Trying alternative loading strategy...")
51        try:
52            # Try 8-bit quantization if 4-bit fails
53            model = AutoModelForCausalLM.from_pretrained(
54                model_name,

```

```

49         device_map="auto",
50         load_in_8bit=True,
51         torch_dtype=torch.float16,
52         offload_folder="offload_folder"
53     )
54     except Exception as e2:
55         print(f"8-bit loading also failed: {e2}")
56         print("Trying with minimal configuration...")
57         # Last resort - try loading with minimal settings
58         model = AutoModelForCausalLM.from_pretrained(
59             model_name,
60             device_map="auto",
61             torch_dtype=torch.float16,
62             low_cpu_mem_usage=True
63         )
64
65     return model, tokenizer

1 def process_data_for_training(data, tokenizer, model_name):
2     """Process dataset with improved chunking for memory
3     efficiency"""
4     print("Processing dataset for training...")
5
6     # Improved field detection logic for different dataset
7     # structures
8     sample = data[0]
9     if 'prompt' in sample and 'response' in sample:
10         question_key, answer_key = 'prompt', 'response'
11     elif 'question' in sample and 'answer' in sample:
12         question_key, answer_key = 'question', 'answer'
13     elif 'input' in sample and 'output' in sample:
14         question_key, answer_key = 'input', 'output'
15     else:
16         # Make an educated guess based on available fields
17         keys = list(sample.keys())
18         str_keys = [k for k in sample.keys() if isinstance(sample[k], str)]
19         if len(str_keys) >= 2:
20             question_key, answer_key = str_keys[0], str_keys[1]
21             print(f"Using inferred fields - Question: '{question_key}', Answer: '{answer_key}'")
22         else:
23             raise ValueError(f"Cannot identify question and answer fields. Keys: {keys}")
24
25     print(f"Using fields - Question: '{question_key}', Answer: '{answer_key}'")
26
27     # Process the data in smaller chunks for better memory
28     # management
29     chunk_size = 100
30     processed_data = []
31
32     for chunk_start in range(0, len(data), chunk_size):
33         chunk_end = min(chunk_start + chunk_size, len(data))
34         print(f"Processing chunk {chunk_start} to {chunk_end-1}...")
35
36         for i in range(chunk_start, chunk_end):

```

```

34         try:
35             item = data[i]
36             question = item[question_key]
37             answer = item[answer_key]
38
39             # Skip invalid entries
40             if not isinstance(question, str) or not isinstance
               (answer, str):
41                 print(f"Skipping item {i}: Invalid data types"
                       )
42                 continue
43
44             # Enhanced prompt engineering with explicit
               reasoning instructions
45             formatted_question = f"""Solve this math problem
               by breaking it down into small, logical steps.
46
47             Problem: {question}
48
49             Follow these steps:
50             1. Understand what the problem is asking for
51             2. Identify the key variables and relationships
52             3. Plan your solution approach step-by-step
53             4. Execute each calculation carefully
54             5. Verify your answer is reasonable
55             6. State the final answer after ####
56
57             Your solution: """
58
59             # Ensure answers end with the #### marker if not
               already present
60             if '####' not in answer:
61                 # Try to find the final numerical answer
62                 numbers = re.findall(r"([-+]?[d*\.]?[d+]",
               answer)
63                 if numbers:
64                     final_number = numbers[-1].strip()
65                     # Only add #### if we're not already at
               the end of the answer
66                     if not answer.strip().endswith(
               final_number):
67                         answer = answer.strip() + f"\n\n#### {
               final_number}"
68
69             # Format for model input
70             if "mistral" in model_name.lower():
71                 formatted_text = f"<s>[INST] {
               formatted_question} [/INST] {answer}</s>"
72             else:
73                 formatted_text = f"<|im_start|>user\n{
               formatted_question}<|im_end|>\n<|im_start
               |>assistant\n{answer}<|im_end|>"
74
75             # Tokenize with longer context window
76             tokenized = tokenizer(formatted_text, truncation=
               True, max_length=2048)
77
78             processed_data.append({
79                 "input_ids": tokenized["input_ids"],

```



```

80         "attention_mask": tokenized["attention_mask"],
81     })
82
83     # Show sample for debugging
84     if i == chunk_start:
85         print(f"\nSample processed data (item {i}):")
86         print(f"Original question: {question[:100]}...
87             " if len(question) > 100 else f"Original
88             question: {question}")
89         print(f"Original answer: {answer[:100]}..." if
90             len(answer) > 100 else f"Original answer:
91             {answer}")
92         print(f"Tokenized length: {len(tokenized['
93             input_ids'])}")
94
95     except Exception as e:
96         print(f"Error processing item {i}: {e}")
97         continue
98
99     # Clear memory after each chunk
100    clear_memory()
101    print(f"Processed {len(processed_data)} examples so far")
102
103    # Create dataset from processed data
104    dataset = Dataset.from_dict({
105        "input_ids": [item["input_ids"] for item in processed_data
106            ],
107        "attention_mask": [item["attention_mask"] for item in
108            processed_data]
109    })
110
111    return dataset

```

```

1  # Enhanced LoRA configuration for Mistral model
2  peft_config = LoraConfig(
3      r=32, # Increased from 8 to 32 for more
4          capacity
5      lora_alpha=64, # Increased from 16 to 64 for stronger
6          updates
7      lora_dropout=0.1, # Increased dropout for better
8          generalization
9      bias="none",
10     task_type="CAUSALLM",
11     target_modules=[
12         # Expanded target modules for more comprehensive
13         # adaptation
14         "q-proj", "k-proj", "v-proj", "o-proj",
15         "gate-proj", "up-proj", "down-proj",
16     ]
17 )
18
19 # Improved training arguments
20 training_args = TrainingArguments(
21     output_dir=output_dir,
22     num_train_epochs=2, # Increased from 1 to 2
23     epochs
24     per_device_train_batch_size=2, # Reduced batch size for
25     memory efficiency

```

```

20     gradient_accumulation_steps=8,      # Increased for larger
        effective_batch
21     learning_rate=1e-4,                # Slightly lower learning
        rate
22     weight_decay=0.05,                 # Increased weight decay
        for regularization
23     warmup_ratio=0.05,                 # Longer warmup phase
24     max_grad_norm=0.5,                 # Increased for stability
25     logging_steps=10,
26     save_strategy="epoch",             # Save checkpoints each
        epoch
27     save_total_limit=3,                 # Keep top 3 checkpoints
28     optim="paged_adamw_32bit",          # Memory-efficient optimizer
29     fp16=True,                         # Mixed precision
30     gradient_checkpointing=True,        # Memory efficiency
31     lr_scheduler_type="cosine",         # Better scheduler
32     report_to="none",                  # Disable wandb/tensorboard
        to save memory
33 )

```

3.3 TINYLLAMA ENHANCED IMPLEMENTATION

For the smaller TinyLlama model, we focused on advanced fine-tuning techniques to maximize performance:

```

1 def format_prompt_for_tinyllama(question):
2     """Format prompt with improved instructions for mathematical
        reasoning"""
3     return f"""<|system|>
4     You are a highly intelligent math assistant that excels at solving
        complex math problems
5     step by step using chain-of-thought reasoning. You always show
        your work clearly,
6     explaining each step of your calculation, and double-check your
        final answer.
7     <|user|>
8     Please solve this math problem by breaking it down into clearly
        defined steps.
9     Show all your work and calculations, and make sure to verify your
        answer.
10
11     Problem: {question}
12
13     First understand the problem, identify what is being asked, plan
        your solution approach,
14     and then solve it step-by-step. After reaching your final answer,
        include it at the end
15     after '####'.
16     <|assistant|>"""

1 # Enhanced LoRA configuration for TinyLlama
2 peft_config = LoraConfig(
3     r=32,                                # Increased from 16 to 32 for more
        capacity
4     lora_alpha=64,                       # Increased from 32 to 64 for stronger
        updates
5     lora_dropout=0.1,                   # Increased dropout for better
        generalization
6     bias="none",

```

```

7     task_type="CAUSALLM",
8     target_modules=[
9         "q-proj", "v-proj", "k-proj", "o-proj",      # Attention
            modules
10        "gate-proj", "up-proj", "down-proj",      # MLP modules
11        "W_pack"                                   # Special
            module for TinyLlama
12    ]
13 )
14
15 # Extended training arguments with longer training
16 training_args = TrainingArguments(
17     output_dir=output_dir,
18     num_train_epochs=3,                          # Increased to 3 epochs
19     per_device_train_batch_size=8,                # Increased batch size for
            smaller model
20     gradient_accumulation_steps=4,                # Increased for larger
            effective batch size
21     learning_rate=1e-4,                          # Slightly lower learning
            rate for stability
22     weight_decay=0.05,                           # Increased weight decay
            for regularization
23     warmup_ratio=0.05,                           # Slightly increased warmup
24     max_grad_norm=0.5,                           # Increased for stability
25     logging_steps=10,
26     save_strategy="epoch",                       # Save per epoch
27     save_total_limit=3,                          # Keep top 3 checkpoints
28     optim="paged_adamw_32bit",                   # Memory-efficient optimizer
29     fp16=True,                                    # Mixed precision
30     gradient_checkpointing=True,                  # Memory efficiency
31     lr_scheduler_type="cosine",                   # Better scheduler for math
            tasks
32     report_to="none",                            # Disable wandb/tensorboard
            to save memory
33 )

```

3.4 KEY IMPLEMENTATION DIFFERENCES

The two implementations differed in several key aspects, reflecting our approach to optimizing for different model sizes:

Component	Mistral-7B Implementation	TinyLlama Enhanced
Model loading	Progressive fallback strategy	Standard loading with 4-bit quantization
Memory optimization	Aggressive CPU offloading	Standard memory optimization
LoRA configuration	Focus on attention layers	Broader parameter coverage
Training duration	2 epochs	3 epochs
Batch size	Smaller (2 per device)	Larger (8 per device)
Prompt engineering	Structured reasoning steps	More detailed reasoning guidance

Table 1: Key implementation differences between model pipelines

These differences highlight our approach to adapting the fine-tuning process based on model size and memory constraints, with more aggressive memory optimization for the larger model and more comprehensive parameter adaptation for the smaller model.

4 EXPERIMENTS

Our experiments with both models revealed significant insights about the relationship between model size, computational resources, and fine-tuning outcomes for mathematical reasoning tasks.

4.1 EVALUATION SETUP

Both models were evaluated using a consistent approach:

Evaluation Dataset:

- 50 examples from the GSM8K benchmark test set
- Problems spanning various mathematical concepts and difficulty levels
- Consistent extraction of ground truth answers for fair comparison

Prompting Strategy:

- Model-specific prompt templates that encourage step-by-step reasoning
- Consistent instruction to provide the final answer after "" markers
- Temperature setting of 0.1-0.2 to balance creativity and determinism

Evaluation Metric:

- Answer accuracy: Exact match between extracted model answer and ground truth
- Normalization of answers to account for formatting differences
- Analysis of solution approaches for correct and incorrect responses

Computational Environment:

- Google Colab with T4 GPU (approximately 16GB VRAM)
- Limited execution time (maximum 12 hours per session)
- Python 3.10 with PyTorch 2.0 and Transformers 4.31.0

4.2 QUANTITATIVE RESULTS

Our evaluation revealed significant performance differences between the models both before and after fine-tuning:

Model	Accuracy	Correct/Total
Base Mistral-7B	0.46	23/50
Fine-tuned Mistral-7B	0.54	27/50
Base TinyLlama-1.1B	0.00	0/50
Fine-tuned TinyLlama-1.1B	0.40	20/50

Table 2: Model performance on GSM8K evaluation set

Training Loss Progression:

- **Mistral-7B:** Training loss dropped from 6.07 to 0.43 by the end of 2 epochs
- **Mistral-7B:** Training loss dropped from 6.07 to 0.38 by the end of 2 epochs
- **TinyLlama-1.1B (Enhanced):** Training loss dropped from 1.27 to 0.39 by the end of 3 epochs

Resource Utilization and Limitations:

- **Mistral-7B:**

- Successfully completed both training and evaluation with enhanced memory optimization techniques
- Training time: approximately 8 hours for 2 epochs
- Inference time: 8.7s/example
- Peak memory usage: 14.2GB of 16GB available VRAM
- **TinyLlama-1.1B (Enhanced):**
 - Complete pipeline executed successfully with fewer memory constraints
 - Training process took approximately 9 hours for 3 epochs
 - Inference time: 3.2s/example
 - Peak memory usage: 8.7GB of 16GB available VRAM

4.3 RESULT ANALYSIS

Several critical insights emerge from our experimental results:

1. Base Model Capability Gap: The stark difference in base model performance (46% vs. 0% accuracy) highlights the dramatic capability gap between 7B and 1.1B parameter models for mathematical reasoning. This substantial disparity suggests that model size plays a crucial role in developing the foundational reasoning capabilities necessary for mathematical problem-solving.

The 0% accuracy of the base TinyLlama model is particularly striking. Detailed inspection of its outputs revealed that while the model could generate text that superficially resembled mathematical solutions with appropriate formatting, it consistently failed to execute correct calculations or arrive at valid answers. The model would often:

- Set up the problem correctly but make fundamental calculation errors
- Lose track of variables midway through multi-step problems
- Generate mathematically inconsistent intermediate steps
- Produce final answers with no logical connection to the solution steps

This suggests that mathematical reasoning may require a minimum model capacity threshold for baseline competence.

2. Advanced Fine-tuning Effectiveness: Our enhanced fine-tuning approach for TinyLlama produced remarkable improvements, boosting accuracy from 0% to 40%. This dramatic 40 percentage point increase demonstrates that even small models can develop specialized capabilities with the right training methodology. The success of our approach can be attributed to several key enhancements:

- **Multi-source training data:** Combining examples from PrimeIntellect/NuminaMath, GSM8K, and reasoning-machines/gsm-hard provided more diverse mathematical patterns and reasoning approaches
- **Improved prompt engineering:** Explicit instructions for step-by-step reasoning helped guide the model toward systematic problem-solving
- **Enhanced LoRA configuration:** Higher rank (32) and alpha (64) values enabled more expressive parameter adaptations
- **Extended training duration:** Three full epochs allowed for more complete learning
- **Optimized learning rate scheduling:** Cosine scheduling with warm-up provided better convergence

The Mistral-7B model also showed improvement, increasing from 46% to 54% accuracy. While this 8 percentage point gain is smaller in absolute terms than TinyLlama's improvement, it represents a relative improvement of 17.4% and demonstrates that our approach benefits larger models as well. The lower relative improvement suggests that larger models may already be operating closer to their optimal performance for this task.

3. Memory Optimization Impact: Our experiments clearly demonstrate the critical role of memory optimization in enabling fine-tuning of larger models:

- The progressive fallback strategy for model loading proved essential for working with the 7B parameter model within T4 GPU constraints
- Chunked data processing enabled handling larger datasets without memory overflow
- Gradient checkpointing and 4-bit quantization significantly reduced memory requirements with minimal performance impact
- CPU offloading for state dictionaries allowed working with larger models than would otherwise be possible

These optimizations made it possible to not only train but also evaluate the larger model within consumer-grade GPU constraints, challenging the notion that high-end computational resources are always necessary.

4. Prompt Engineering Importance: Detailed analysis of model outputs revealed that improved prompt engineering had a substantial impact on performance:

- Models prompted with explicit step-by-step reasoning instructions produced more structured solutions
- The inclusion of specific verification steps improved calculation accuracy
- Consistent formatting of answers (using markers) improved answer extraction reliability
- Enhanced prompts particularly benefited the smaller model, providing more explicit guidance

This finding suggests that careful prompt design can partially compensate for model size limitations, especially for tasks with well-defined solution structures like mathematical reasoning.

5 CONCLUSION

Our research into advanced fine-tuning strategies for mathematical reasoning in resource-constrained environments yields several important insights:

1. Optimized Memory Management as an Enabler: Our implementation of progressive fallback strategies and advanced memory optimization techniques for the Mistral-7B model demonstrates that larger models can be successfully fine-tuned and evaluated even within the constraints of consumer-grade GPUs like the T4. These techniques included 4-bit quantization, gradient checkpointing, CPU offloading, and dynamic memory allocation based on available resources. This finding challenges the notion that high-end computational resources are always necessary for working with larger language models.

2. Model Size vs Training Methodology Trade-off: While the 7B-parameter Mistral model achieved the highest overall performance (54% accuracy after fine-tuning), our enhanced fine-tuning approach enabled the much smaller 1.1B-parameter TinyLlama to achieve 40% accuracy from a starting point of 0%. The 8 percentage point improvement for Mistral compared to the 40 percentage point improvement for TinyLlama suggests that smaller models may benefit disproportionately from advanced fine-tuning techniques. This demonstrates that methodological improvements can substantially compensate for parameter count limitations, challenging the notion that larger models are always necessary for complex reasoning tasks.

3. Data Diversity as a Critical Factor: The enhanced implementation leveraged multiple data sources, which proved crucial for developing robust mathematical reasoning. This multi-source approach provided exposure to diverse problem types and solution strategies, enabling more generalizable learning across both models.

4. Prompt Engineering as a Performance Multiplier: Our implementation incorporated detailed prompts with explicit instructions for step-by-step reasoning, which significantly improved both models' ability to structure and execute mathematical solutions. This finding suggests that all models, regardless of size, benefit from clear guidance in both training and inference contexts.

5. Parameter-Efficient Fine-Tuning Effectiveness: Our success with enhanced LoRA configurations (increased rank and alpha values) demonstrates that parameter-efficient methods can be

remarkably effective when properly optimized. By concentrating adaptation on key model components and using higher-capacity LoRA parameters, we achieved substantial improvements without increasing the overall parameter count or memory requirements.

These findings have important implications for researchers and practitioners working with LLMs under resource constraints:

- When computational resources are limited, focusing on memory optimization techniques, data quality, prompt engineering, and fine-tuning hyperparameters can yield substantial improvements across models of different sizes.
- For specialized tasks like mathematical reasoning, a carefully fine-tuned smaller model may perform comparably to a larger base model, offering a practical alternative when computational resources are constrained.
- Multi-source training data provides exposure to more diverse problem-solving approaches, which appears especially beneficial for developing robust reasoning capabilities.
- Enhanced LoRA configurations with higher rank and alpha values can significantly improve adaptation capacity without substantially increasing memory requirements.
- Progressive fallback strategies for model loading enable more reliable experimentation with larger models in constrained environments.

In conclusion, our research demonstrates that advanced fine-tuning strategies can dramatically improve the performance of language models of various sizes on complex reasoning tasks. While computational resources remain an important factor in LLM fine-tuning, our work shows that methodological innovation can substantially mitigate resource limitations, potentially democratizing access to specialized AI capabilities.

ACKNOWLEDGMENT

This is the Assignment 3 for CSC6052 / MDS5110 / CSC5051, see details in <https://nlp-course-cuhksz.github.io/>.

REFERENCES