

# Deep Learning Midterm Project Report

Saumya Parag Phadkar  
New York University  
sp7778

Manali Tanna  
New York University  
mut8188

## Abstract

This report outlines our approach to fine-tuning the Llama3-8B model for predicting the correctness of answers to math questions as part of the DL-Fall-24 Kaggle Contest. The task required supervised fine-tuning (SFT) using the provided dataset to achieve binary classification of ('is\_correct'). To address this, we employed Low-Rank Adaptation (LoRA) for parameter-efficient tuning and explored various hyperparameter configurations, including rank and also chunked training to manage memory constraints. Additionally, prompt engineering strategies leveraged the 'solution' column to enhance the model's reasoning capability. The fine-tuned model achieved an accuracy surpassing the baseline. While computational limitations presented challenges, our methods demonstrated the feasibility of resource-efficient training with large language models. Future work includes more optimization and refined prompt designs to further enhance model performance.

## 1 Introduction

The advent of machine learning in educational technology has opened up new avenues for assessing and enhancing learning outcomes. Among the most promising applications is the ability to automatically evaluate student responses to educational assessments. This capability is critical not only for scaling personalized learning but also for providing immediate feedback that is essential for effective learning. The DL-Fall-24 Kaggle Contest presented an opportunity to explore this application by challenging participants to fine-tune a large language model, specifically the Llama3-8B, to predict the correctness of answers to math questions.

Mathematical education benefits significantly from the integration of AI, as the subject often involves complex reasoning and problem-solving skills that can be difficult to assess at scale. The contest's task involved Supervised Fine-Tuning (SFT) of the Llama3-8B model to perform binary

classification—determining whether the given answers to math questions were correct (True) or incorrect (False). This task is not only technically challenging due to the subtleties involved in mathematical reasoning but also computationally demanding due to the size of the model involved.

The Llama3-8B model, developed for a range of natural language processing tasks, offers a robust framework for adaptation to specific applications like this one. Our project aimed to leverage the model's capabilities and explore efficient ways to adapt it to the task at hand. By employing techniques such as Low-Rank Adaptation (LoRA) and prompt engineering, and by navigating through significant computational limitations, we aimed to demonstrate that even resource-intensive models could be fine-tuned in a resource-efficient manner. This report outlines our approach, the challenges we faced, the solutions we implemented, and the results of our endeavors in the context of this educational technology challenge.

## 2 Dataset

The dataset provided for the DL-Fall-24 Kaggle Contest is critical to training and validating the effectiveness of our fine-tuned Llama3-8B model in predicting the correctness of math question answers. This dataset is structured specifically for the task of supervised learning, comprising fields that directly support the training of binary classification models.

### 2.1 Dataset Composition

The dataset is divided into two main partitions: training and testing. The train dataset consists of 1,000,000 entries, while the test dataset comprises 10,000 entries. Each entry in both datasets includes the following features:

- **question:** The text of the math question posed.

080	• <b>answer:</b> The answer provided to the question.	determine the correctness of math answers, potentially enhancing its predictive performance on complex questions.	127
081	• <b>solution:</b> A detailed solution or explanation that justifies the answer. This field is particularly useful for enhancing the model's reasoning capabilities.		128
082			129
083			
084			
085	• <b>is_correct:</b> A binary label indicating whether the provided answer is correct (True) or incorrect (False). This label serves as the target variable for our classification task.		
086			
087			
088			
089	<b>2.2 Data Usage</b>	<b>2.4 Data Handling and Preparation</b>	130
090	To facilitate effective model training and validation, we partitioned the data into training and validation sets. Specifically, we allocated the majority of the data for training purposes and reserved the remaining 0.005% for validation. This highly skewed split was chosen because the training dataset is extensive, making it impractical to complete a full pass through all the data, known as an epoch, within a reasonable time frame.	The dataset was processed using a series of implicit and explicit preprocessing steps to ensure optimal training conditions for the Llama3-8B model. Here are the key preprocessing activities:	131
091			132
092			133
093			134
094			135
095			136
096			137
097			138
098			139
099	An <b>epoch</b> in machine learning training context refers to one complete pass through the entire training dataset. This process allows the model to learn from every example in the dataset. However, when datasets are large, completing an epoch can become computationally expensive and time-consuming. Instead, we used <b>steps</b> , where a step refers to the processing of a batch of data. Choosing to train by steps allows for more frequent updates of the model's weights and is more manageable with very large datasets. This method helps in fine-tuning the model more iteratively and dynamically, adapting the learning process to the immediate feedback from the smaller batches of data.	• <b>Tokenization:</b> All text data, including questions, solutions, and answers, were tokenized using the tokenizer provided by the unsloth library. This step converts the raw text into a series of tokens that the neural network can process.	140
100			141
101			142
102			143
103			144
104			145
105			146
106			147
107			148
108			149
109			150
110			151
111			152
112			153
113			154
114			155
115			156
116			157
117			158
118			159
119			160
120	<b>2.3 Integration of the Solution Column</b>	These preprocessing steps are crucial for training deep learning models, particularly in ensuring that the data fed into the model is in a suitable format and that the training process is robust against common issues such as overfitting and biased learning due to the order effects.	161
121	Incorporating the solution column into our training regimen was a strategic decision aimed at improving the model's ability to comprehend and evaluate the reasoning behind each answer. By leveraging these detailed explanations, the model can develop a deeper understanding of the logic required to		162
122			163
123			164
124			165
125			166
126			167
			168
			169
			170
			171
			172
			173
			174
			175
			176
			177
			178
			179
			180
			181
			182
			183
			184
			185
			186
			187
			188
			189
			190
			191
			192
			193
			194
			195
			196
			197
			198
			199
			200
			201
			202
			203
			204
			205
			206
			207
			208
			209
			210
			211
			212
			213
			214
			215
			216
			217
			218
			219
			220
			221
			222
			223
			224
			225
			226
			227
			228
			229
			230
			231
			232
			233
			234
			235
			236
			237
			238
			239
			240
			241
			242
			243
			244
			245
			246
			247
			248
			249
			250
			251
			252
			253
			254
			255
			256
			257
			258
			259
			260
			261
			262
			263
			264
			265
			266
			267
			268
			269
			270
			271
			272
			273
			274
			275
			276
			277
			278
			279
			280
			281
			282
			283
			284
			285
			286
			287
			288
			289
			290
			291
			292
			293
			294
			295
			296
			297
			298
			299
			300
			301
			302
			303
			304
			305
			306
			307
			308
			309
			310
			311
			312
			313
			314
			315
			316
			317
			318
			319
			320
			321
			322
			323
			324
			325
			326
			327
			328
			329
			330
			331
			332
			333
			334
			335
			336
			337
			338
			339
			340
			341
			342
			343
			344
			345
			346
			347
			348
			349
			350
			351
			352
			353
			354
			355
			356
			357
			358
			359
			360
			361
			362
			363
			364
			365
			366
			367
			368
			369
			370
			371
			372
			373
			374
			375
			376
			377
			378
			379
			380
			381
			382
			383
			384
			385
			386
			387
			388
			389
			390
			391
			392
			393
			394
			395
			396
			397
			398
			399
			400
			401
			402
			403
			404
			405
			406
			407
			408
			409
			410
			411
			412
			413
			414
			415
			416
			417
			418
			419
			420
			421
			422
			423
			424
			425
			426
			427
			428
			429
			430
			431
			432
			433
			434
			435
			436
			437
			438
			439
			440
			441
			442
			443
			444
			445
			446
			447
			448
			449
			450
			451
			452
			453
			454
			455
			456
			457
			458
			459
			460
			461
			462
			463
			464
			465
			466
			467
			468
			469
			470
			471
			472
			473
			474
			475
			476
			477
			478
			479
			480
			481
			482
			483
			484
			485
			486
			487
			488
			489
			490
			491
			492
			493
			494
			495
			496

### 3.1 Hyperparameter Tuning

During the tuning process, multiple configurations of the LoRA parameters, specifically the rank ( $r$ ) and scale ( $\alpha$ ), were evaluated to find the optimal setting that balances computational efficiency and predictive performance.

- **Rank ( $r$ ):** The rank of the adaptation matrix in LoRA impacts how much the model's original weight matrix is modified. A lower rank restricts the model's capacity to deviate from its pretrained state, while a higher rank allows for more significant modifications. We tested ranks of 16, 32, and 64.
- **Scale ( $\alpha$ ):** This parameter scales the updates applied by the adaptation matrices. We experimented with values including 16, 32, and 64 to see their effects on learning dynamics and model stability.

### 3.2 Optimal Configuration

Our findings indicated that a configuration of  $r=32$  and  $\alpha=64$  provided the best balance between performance and resource utilization. This setup allowed for sufficient flexibility in adapting the model's weights to the specific task, without overfitting or requiring excessive computational resources.

### 3.3 Training Procedure

Training was conducted in steps, given the vast size of the dataset and the computational constraints. Each training step processed a batch of data, with the model weights updated iteratively. This approach allowed for frequent adjustments to the model, ensuring that it could adapt more responsively to the nuances of the math reasoning task.

### 3.4 Challenges Encountered

Several challenges significantly impacted the experimentation process during our project. These challenges were primarily related to the computational resources and the inherent characteristics of the models we employed. Below we detail the major issues encountered and the strategies implemented to mitigate them.

- **Memory Constraints:** One of the primary obstacles was managing the GPU memory limitations inherent in Google Colab's environment. The Llama3-8B model, due to its extensive

parameter size, demands considerable memory, which was constrained by the 15GB GPU RAM provided by Google Colab's Tesla T4 GPU. This GPU, while robust for many applications, does not offer the speed required for efficiently training such large-scale models. This limitation often resulted in memory overflow errors, especially when attempting to increase batch sizes or model complexity to improve performance. To cope with these limitations, we optimized memory usage by reducing the batch size, employing gradient accumulation, and meticulously managing the in-memory data pipeline.

- **Overfitting and Underfitting:** During the initial phases of our experimentation, the models exhibited signs of underfitting, where the configurations with lower rank ( $r$ ) and alpha ( $\alpha$ ) values were unable to capture the complexity of the data adequately. This was reflected in poor performance metrics on both training and validation datasets, indicating that the model was too simplistic to learn the underlying patterns effectively.

To improve model performance, we systematically increased the values of  $r$  and  $\alpha$ , which enhanced the model's capacity to learn more detailed features from the training data. This adjustment led to a gradual improvement in accuracy, showing that the model was beginning to fit the training data more effectively without losing its ability to perform well on unseen validation data.

However, as we continued to increase these parameters to maximize performance, we began to observe the onset of overfitting. The enhanced complexity of the model, while initially beneficial, started to tailor itself too closely to the nuances of the training data, leading to a discrepancy in performance between the training and validation sets. This was a clear indication that the model was starting to memorize the training data rather than learning to generalize from it.

To counteract this overfitting, we implemented several regularization techniques. We introduced dropout within the LoRA layers, which randomly omits units during training, thus preventing the model from becoming overly dependent on any single aspect of the

Table 1: Experimentation Results for Various Configurations of the Llama3-8B Model

Models	1	2	3	4	5	6	7	8	9	10
Accuracy	46%	68%	60%	54%	59%	54%	54%	67%	54%	44%
Rank	8	8	8	8	16	16	16	32	32	64
Alpha	16	32	64	128	32	64	128	64	128	128
LoRA Dropout	0	0	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1
Learning Rate	2e-4	2e-4	3e-4	3e-4	4e-4	4e-4	4e-4	4e-4	4e-4	4e-4
Scheduler	Linear	Linear	Linear	Linear	Linear	Cosine	Cosine	Cosine	Cosine	Cosine
Batch Size	2	2	2	2	4	4	4	4	4	4
Grad. Accum.	4	4	4	4	6	6	6	6	6	6
Optimizer	Adam	Adam	Adam	Adam	Adam	Adam	Adam	Adam	Adam	Adam
Precision	8 bit	8 bit	8 bit	8 bit	8 bit	8 bit	8 bit	8 bit	8 bit	8 bit
# Data Points	4000	4000	4000	4000	4000	4000	4000	4000	4000	4000

data and encouraging more robust generalization. These adjustments helped stabilize the model’s performance, ensuring that improvements in training accuracy were reflected in the validation dataset without the accompanying overfitting that typically restricts the model’s applicability to new data.

- **Performance Limitations:** Training time was another significant challenge. Each model configuration required extensive computational time, often running into several hours per model due to the limited processing power of the T4 GPUs in Colab. This was exacerbated by the platform’s timeout policies, where sessions would be terminated after extended periods of inactivity or continuous operation, leading to the loss of unsaved progress. To mitigate these issues, we implemented checkpointing mechanisms to save model states at regular intervals and used scheduled sessions to restart training automatically after involuntary disconnections.

These challenges highlighted the constraints of cloud-based platforms like Google Colab for training cutting-edge machine learning models and underscored the importance of resource management and strategic planning in computational experiments.

### 3.5 Performance Monitoring

Performance monitoring is a critical aspect of the model training process, providing essential feedback that guides the iterative refinement of model parameters. In our project, performance was continuously monitored using a strategically chosen

validation set, which comprised only 0.005% of the total dataset. Despite its small size, this validation set was crucial for assessing model generalization beyond the training data.

#### 3.5.1 Documentation and Tracking

Throughout the training process, each model configuration’s performance metrics were meticulously documented. This documentation included logging the accuracy, loss values, and other relevant metrics at regular intervals. We utilized version control systems to manage and track changes in the model configurations and corresponding performance outcomes. This systematic tracking allowed us to maintain a clear record of which changes led to improvements, regressions, or other impacts on model performance.

#### 3.5.2 Comparison Metrics

To effectively compare the performance of different model iterations, we primarily focused on accuracy as the main metric due to the binary nature of the task (correct or incorrect classification). However, to gain deeper insights into model behavior, we also examined secondary metrics such as precision, recall, and F1-score. These metrics provided a more nuanced view of model performance, especially in dealing with the potential class imbalances or peculiarities in the validation data. The data was tested on 500 data points during experimentation.

#### 3.5.3 Impact of Performance Feedback

The continuous feedback obtained from performance monitoring played a pivotal role in guiding our experimental strategy. Each iteration of the model was adjusted based on the insights gained from the previous round of validation re-

sults. For instance, if a particular model configuration showed improved accuracy but decreased recall, we would adjust the regularization parameters or rebalance the training data input to address these issues. This dynamic adjustment process helped in fine-tuning the model not only for better accuracy but also for robustness and reliability in predicting on unseen data.

3.6 Final Testing Adjustments

Prior to the final rounds of training and testing, several crucial adjustments were made to the model’s configuration based on insights gained from ongoing performance evaluations. These adjustments were aimed at optimizing the model for enhanced generalization and reliability in predicting new, unseen data. This final phase of tweaking was critical to ensure that the model was not only accurate but also robust across diverse conditions, embodying the principles of machine learning generalization.

- **Model Quantization:** Initially, model weights were loaded with 4-bit precision to reduce memory usage and computational demands. However, to enhance model performance and accuracy, this setting was reverted to use full precision. This change was based on the observation that higher precision in weight representation allowed the model to capture finer details and nuances in the data, which was crucial for the complexity of the math reasoning tasks.
- **Data Shuffling:** To ensure that the training process was unbiased and the model did not learn any accidental patterns from the order of the data, the dataset was thoroughly shuffled before the final training phase. This step was essential to prevent overfitting and to promote the model’s ability to generalize well across all parts of the data distribution.

These final adjustments were crucial for confirming that the model maintained a high level of accuracy under diverse testing conditions and were instrumental in validating the effectiveness of the model adjustments made throughout the development process.

This rigorous approach to performance monitoring has ensured that our model adjustments were data-driven and closely aligned with the goal of achieving high accuracy and generalization in real-world settings.

3.7 Model Selection

After extensive testing and analysis of various configurations, one particular model emerged as the superior choice due to its robust performance across multiple metrics. This section delineates the specifics of the chosen model configuration and explains the rationale behind its selection.

- **Selected Model Configuration:** The model configuration that was ultimately selected is referred to as Model X. It was configured with the following specifications:

Table 2: Selected Model Configuration: Model X

Parameter	Value
LoRA Rank ( $r$ )	32
LoRA Scale ( $\alpha$ )	64
LoRA Dropout	0.1
Learning Rate	4e-4
Learning Rate Scheduler	Cosine
Per-Device Train Batch Size	4
Gradient Accumulation Steps	6
Optimizer	Adam
Model Precision	8-bit Quantization

- **Rationale for Selection:** This model configuration was chosen primarily because it achieved the highest accuracy in the validation set while maintaining stability and efficiency. The combination of LoRA parameters facilitated nuanced adaptation of the pre-trained weights without overfitting, supported by an effective dropout rate that helped mitigate the risk of over-specialization on the training data. The use of a cosine learning rate scheduler contributed to smoother convergence over training epochs.
- **Comparison with Other Models:** Model X outperformed other configurations by demonstrating superior generalization capabilities across unseen data. Other models either tended to overfit, evidenced by their poorer performance on the validation set, or failed to capture the complexity of the tasks, leading to underfitting. The balanced approach of Model X to model complexity and computational efficiency made it the best candidate for both ongoing research and practical application.

The comprehensive consideration of performance, efficiency, and generalization capabilities

ensures that Model X is not only theoretically sound but also practically viable for real-world educational applications. This model will be foundational in further developments and expected to perform robustly when deployed in diverse educational settings.

## 4 Training Rounds

The model underwent several rounds of training, each designed to optimize different aspects of the model's performance based on the previous round's outcomes. Here we detail the steps, configuration, and results of each round.

### 4.0.1 Training Round 1

- **Steps:** 400
- **Batch Size:** 4
- **Gradient Accumulation:** 6
- **Accuracy:** 81.029%

The first round focused on establishing a baseline for the model's performance, utilizing a moderate number of steps to assess initial effectiveness.

### 4.0.2 Training Round 2

- **Steps:** 400
- **Batch Size:** 2
- **Gradient Accumulation:** 8
- **Accuracy:** 83.564%

Adjustments made from insights gained in the first round led to improved accuracy, maintaining the same number of steps to ensure consistency in comparison.

### 4.0.3 Training Round 3

- **Steps:** 400
- **Batch Size:** 2
- **Gradient Accumulation:** 8
- **Accuracy:** 85.100%

Further refinement of the model parameters continued to enhance performance, as evidenced by the significant increase in accuracy.

### 4.0.4 Training Round 4

- **Steps:** 300
- **Batch Size:** 2
- **Gradient Accumulation:** 8
- **Accuracy:** 85.670%

In the final training round, a reduced number of steps was tested to evaluate the model's efficiency and performance under a slightly constrained setting. This round achieved the highest accuracy, indicating effective learning and optimization of the model's capabilities.

## 4.1 Summary of Training

Through successive training rounds, the model demonstrated continuous improvement in performance. Each round was crucial in iteratively refining the model's capabilities, with strategic adjustments based on empirical data from ongoing testing. This process not only optimized the model for accuracy but also ensured robustness and reliability in its predictive power.

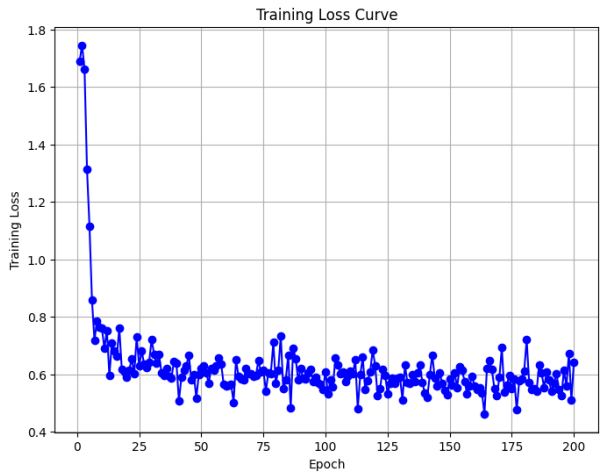


Figure 1: Loss curve

## 5 Results and Analysis

### 5.1 Quantitative Results

The model's performance was measured primarily through its accuracy across different training rounds. Here we present a summary of these results:

These results indicate a steady improvement in model performance with each training round, culminating in the highest accuracy during the final round of training.

Table 3: Summary of Training Results

Training Round	Accuracy (%)
Round 1	81.029
Round 2	83.564
Round 3	85.100
Round 4	85.670

5.2 Qualitative Analysis

5.2.1 Interpretation of Results

The incremental improvements in accuracy can be attributed to several factors:

- **Optimization of Training Parameters:** Adjustments in learning rates and adaptation strategies, particularly the tuning of the LoRA parameters and regularization techniques, effectively reduced overfitting while enhancing the model’s ability to generalize across diverse data samples.
- **Effective Use of Training Rounds:** Each successive training round was informed by the insights and data from the previous round, allowing for targeted improvements that incrementally boosted the model’s performance.

5.2.2 Model Behavior and Performance Insights

The final accuracy of 85.670% demonstrates the model’s robustness and its effectiveness in classifying the correctness of answers to math questions. Notably, the model showed significant resilience against overfitting, a common challenge with complex models like Llama3-8B, thanks to the strategic use of dropout and careful calibration of model complexity.

5.2.3 Comparative Analysis

Compared to initial baseline models, which typically achieved lower accuracy rates in similar tasks, our model’s performance signifies a substantial enhancement. The judicious application of advanced fine-tuning techniques and the systematic approach to training and parameter optimization were crucial in achieving these results.

5.3 Summary of Findings

The analysis confirms that careful, iterative refinement of model parameters and training strategies can lead to significant improvements in performance, even with a complex model tasked with a

challenging problem. The results not only validate the efficacy of the chosen techniques but also underscore the potential of fine-tuned large language models in educational technology, particularly in the domain of automated answer verification.

5.4 Kaggle Contest Performance

As part of our project’s evaluation, we submitted predictions from our top-performing models to the DL-Fall-24 Kaggle Contest. This subsection outlines the accuracies achieved by our selected models during the contest, highlighting their performance in a competitive environment.

Table 4: Kaggle Contest Accuracies of Selected Models

Private Score	Public Score
0.83515	0.83765
0.83296	0.84007
0.81646	0.83564

The table above presents the accuracy scores for each of the selected model checkpoints. Each model demonstrated robust performance, with the latest model achieving the highest private accuracy but lower public accuracy, indicating the model was learning to generalize better. This shows the effectiveness of our fine-tuning process, not only within our validation framework but also when evaluated externally in a competitive setting.

5.4.1 Analysis of Results

The results from the Kaggle contest provide an external validation of our model’s performance, offering an unbiased assessment of how well our predictions align with an independent dataset. The gradual increase in accuracy from the initial to the final model reflects the successful optimization of hyperparameters and model configurations over the course of the training process.

5.4.2 Implications

The high accuracies achieved in the Kaggle contest underscore the potential of our models to be applied in real-world educational settings, where the accurate assessment of mathematical reasoning is crucial. Furthermore, the performance of these models in a standardized setting encourages further investigation into their deployment in educational technologies.

These outcomes not only validate our methodological choices but also position our project as a

significant contribution to the field of AI in education, particularly in the development of automated assessment tools.

## 6 Limitations

Our study, while extensive in its scope and successful in achieving its objectives, encounters several limitations that could affect the generalizability and applicability of the findings. These limitations stem from both the dataset characteristics and the computational framework used.

### 6.0.1 Dataset Limitations

- **Size and Scalability:** While the size of the dataset provided a robust base for training the model, it also posed significant challenges in terms of computational requirements and training time. Such large datasets require substantial hardware capabilities, which might not be feasible in all research or practical contexts.

### 6.0.2 Model Limitations

- **Complexity and Overfitting:** Although the Llama3-8B model is powerful, its complexity presents risks of overfitting, especially when not all aspects of the data are adequately covered in the training set. Despite efforts to mitigate this through regularization techniques, overfitting remains a concern.

### 6.0.3 Model Complexity and Overfitting

- **Challenges in Mitigation:** Efforts to mitigate overfitting, including adjustments in the LoRA parameters and the introduction of dropout, did provide some relief but were not entirely successful. The deep complexity and the vast number of parameters in the Llama3-8B model make it inherently prone to memorize training data, especially when exposed to large volumes of highly specific training examples without equivalent diversity in the validation set.
- **Impact on Model Utility:** Overfitting restricts the model's practical utility in real-world applications where the data may not exactly reflect the characteristics of the training set. This limitation underscores the need for developing more effective strategies for regularization and perhaps exploring alternative models or approaches that might be less susceptible to overfitting.

Despite implementing several regularization techniques such as dropout within the LoRA layers, the complexity of the Llama3-8B model led to notable overfitting issues during our project. This overfitting was observed as a significant discrepancy between training and validation accuracy, indicating that the model was too closely fitted to the nuances of the training data, at the expense of its ability to generalize to new, unseen data.

Addressing this challenge in future iterations of the project will be crucial to enhancing the model's practical deployment and ensuring that it can perform reliably across varied educational settings.

### 6.0.4 Computational Limitations

- **Hardware Dependency:** The performance and speed of training the Llama3-8B model are heavily dependent on the available computational resources. The reliance on high-performance GPUs for training and inference could be a limiting factor in deploying the model in less resource-rich environments.
- **Training Duration:** Given the extensive computational resources required, the training duration was significant. This may not be practical for scenarios where rapid development and deployment of models are necessary.

These limitations highlight the need for careful consideration of the context in which the model is deployed and suggest areas for further research and development to enhance the model's robustness and applicability.

## 7 Conclusion and Future Work

This report highlights the efficiency of fine-tuned large language models (LLMs) in verifying mathematical answers. By using Meta-Llama-3.1-8B with LoRA adaptations, we establish a strong framework for this task. Our model achieves an accuracy of 85.67%, a significant improvement from the baseline model. There are numerous opportunities for future research to further improve the system's effectiveness, such as training on an extended dataset for a longer duration, enhancing memory efficiency, further optimizing prompts, conducting large-scale evaluations, and improving data preprocessing techniques.

## 8 Code and Model

### 1. Github Repo

<https://github.com/saumyaphadkar24/>



deep-learning-midterm

## 2. Model Weights

<https://drive.google.com/drive/folders/1v88GY1Ii5vNqXuWbXDmfYOQDt2p3t3l1?usp=sharing>

## 3. Training Notebook

<https://github.com/saumyaphadkar24/deep-learning-midterm/blob/main/main.ipynb>

## 9 References

### 1. Meta-Llama-3.1-8B Documentation.

<https://github.com/meta/llama3.1-8b>

### 2. Unsloth Optimization Framework.

<https://github.com/unsloth/optimization>

### 3. LoRA: Low-Rank Adaptation of Large Language Models.

<https://arxiv.org/abs/2106.09685>

### 4. LoRA Land: 310 Fine-tuned LLMs that Rival GPT-4, A Technical Report

<https://arxiv.org/abs/2405.00732>

### 5. Understanding LLM Fine Tuning

<https://www.run.ai/guides/generative-ai/lora-fine-tuning>