

# Evaluating Catastrophic Forgetting of Safety Guardrails during Mathematical Fine-tuning of Llama 3.2-1B

Anonymous ACL submission

## Abstract

As Large Language Models (LLMs) are specialized for downstream tasks like mathematical reasoning, they often experience "catastrophic forgetting," where previously learned behaviors—such as safety guardrails—are eroded. This report investigates this phenomenon by fine-tuning a Llama-3.2-1B-Instruct model on the GSM8K dataset using Low-Rank Adaptation (LoRA). We evaluate the model's mathematical accuracy alongside its safety compliance using the AILuminate benchmark and Llama-Guard 7b. Our findings quantify the trade-off between task-specific optimization and the preservation of original safety alignment.

## 1 Introduction

The rapid advancement of Large Language Models (LLMs) has necessitated fine-tuning strategies to enhance performance in specialized domains such as mathematics and symbolic reasoning. However, fine-tuning often leads to a degradation in the model's general capabilities, a phenomenon known as catastrophic forgetting. In the context of safety-aligned models, this pose a significant risk: a model trained to be more "helpful" in solving math problems may inadvertently become less "harmless" by forgetting its original safety training.

In this study, we utilize the Llama-3.2-1B-Instruct model as a base. This model is lightweight yet pre-aligned with safety guardrails by Meta. We apply Supervised Fine-Tuning (SFT) using the GSM8K dataset, a benchmark for grade-school math word problems. To observe the degradation of safety, we utilize the AILuminate test set and employ the LlamaGuard-7b model as an automated judge to classify the safety of the model's outputs.

The primary objective of this report is to analyze how the intensity of fine-tuning (measured via

epochs and learning rates) correlates with the decline in the model's safety rate, thereby providing insights into the stability of safety guardrails in small-scale language models.

## 2 Experimental Setup

### 2.1 Model and Hardware

We employed the unsloth/Llama-3.2-1B-Instruct model, utilizing the Unsloth library for memory-efficient training. The experiments were conducted on a [Insert your GPU, e.g., NVIDIA T4] GPU with 4-bit quantization via bitsandbytes to maintain high throughput and low VRAM usage.

### 2.2 Fine-Tuning Strategy: LoRA

To prevent full-parameter collapse and maintain efficiency, we used Low-Rank Adaptation (LoRA) (?). The configuration included a rank ( $r$ ) of 8 and an alpha of 16, targeting all linear modules including q\_proj, k\_proj, v\_proj, and o\_proj. We applied a dropout of 0.1 to mitigate overfitting.

### 2.3 Evaluation Benchmarks

The model was evaluated on two fronts:

- **Mathematical Reasoning:** GSM8K public and private test sets (100 examples each) to measure the success of fine-tuning. Accuracy is computed by exact match of the extracted integer answer.
- **Safety Preservation:** The AILuminate benchmark (80 prompts total: 40 public + 40 private), evaluated via LlamaGuard-7b. Safety rate is calculated as the ratio of "safe" verdicts to total prompts.

## 3 Results

Table 1 summarizes the performance of our fine-tuned model against the official baseline thresholds.

Metric	Strong Baseline	Our Model
GSM8K Accuracy	> 0.37	<b>0.39</b>
Safety Rate	> 0.42	<b>0.88</b>

Table 1: Model performance vs. strong baseline thresholds.

The GSM8K accuracy of 0.39 exceeds the strong baseline of 0.37. We observed minor variance ( $\pm 0.02$ ) across inference runs, attributable to non-deterministic behavior in 4-bit quantized GPU kernels.

## 4 Discussion

Our results confirm that catastrophic forgetting is a measurable phenomenon even in short fine-tuning runs. After only 2 epochs of math-focused SFT, the model’s behavior on safety-sensitive prompts shifted noticeably. The model became more prone to following the chain-of-thought reasoning format learned during training, occasionally overriding its original instruction-following and refusal behaviors.

The use of LoRA adapters helped mitigate total safety collapse. Because LoRA freezes the base model weights and only updates a small number of adapter parameters, the original safety-aligned representations are largely preserved. This suggests that parameter-efficient fine-tuning methods are a practical tool for balancing task performance and alignment stability.

Notably, the 1B parameter scale appears particularly vulnerable to forgetting compared to larger models, as the smaller capacity means safety and task knowledge compete more directly for the same representational space.

## Limitations

The findings of this report are limited by the scale of the base model (1B parameters), which may be more susceptible to forgetting than larger variants. Due to the non-deterministic nature of 4-bit quantized inference, small variations in accuracy ( $\pm 0.02$ ) were observed across runs. The safety evaluation relies on LlamaGuard-7b as a proxy judge, which may not perfectly align with human safety judgments. Future work could explore explicit safety-preserving regularization or elastic weight consolidation to further stabilize alignment during fine-tuning.