

Automated Prompt Optimization on Ultra-Low-Resource Hardware: A Feasibility Study using Llama-3.2-3B

1st Muhammad Zaid

Roll No. 22I-1934

Dept. of Computer Science
FAST-NUCES
Lahore, Pakistan

2nd Abubakar Shahid

Roll No. 22I-1883

Dept. of Computer Science
FAST-NUCES
Lahore, Pakistan

Abstract—The democratization of Artificial Intelligence depends on the ability to run and optimize Large Language Models (LLMs) on consumer-grade hardware. This study reproduces the "Promptomatix" framework under strict resource constraints (Google Colab T4 GPU). While original studies utilized enterprise-grade models (GPT-4), we demonstrate that hardware limitations necessitate the use of highly quantized, lightweight architectures. We successfully operationalized the DSPy optimization pipeline using the Llama-3.2-3B-Instruct model. However, our experiments reveal a "Performance Inversion," where the zero-shot baseline (85.0%) outperformed the few-shot optimized prompt (80.0%). This validates the "Quantization Tax" hypothesis: for highly quantized Small Language Models (SLMs), the noise introduced by automated reasoning traces in the context window outweighs the benefits of few-shot learning, suggesting that simpler prompting strategies are superior for low-resource environments.

Index Terms—Prompt Engineering, DSPy, Quantization, Llama-3, Low-Resource LLMs, Automated Optimization

I. INTRODUCTION

Effective utilization of Large Language Models (LLMs) is critically dependent on prompt engineering. Recent advancements have introduced Automated Prompt Optimization (APO) frameworks, such as Promptomatix [1], which use algorithms like MIPRO and BootstrapFewShot to iteratively refine prompts. These frameworks traditionally rely on high-precision, large-parameter models (e.g., GPT-4, Llama-3-70B) to generate and evaluate reasoning traces.

However, a significant gap exists in applying APO to constrained environments. Students and independent researchers often rely on free-tier infrastructure, such as Google Colab, which offers limited VRAM (16GB). This study explores the feasibility of running a full APO pipeline in such an environment. We detail a reproduction effort that initially faced fatal memory fragmentation with standard 8B models, necessitating a transition to the unsloth-optimized Llama-3.2-3B architecture. We analyze the trade-offs between model size, quantization, and optimization efficacy.

II. METHODOLOGY

A. Hardware Constraints & Model Selection

The experiment was conducted on a single NVIDIA T4 GPU. Initial attempts to load Llama-3-8b-Instruct-bnb-4bit failed due to CUDA OutOfMemory (OOM) errors, as the 8B parameter model requires ~5.7GB VRAM, leaving insufficient overhead for the optimization batches and context caching required by DSPy.

To resolve this, we utilized the unsloth/Llama-3.2-3B-Instruct-bnb-4bit model. This architecture uses 4-bit NF4 quantization to fit within approximately 2.4GB of VRAM. This 60% reduction in memory footprint allowed us to maintain a larger context window for the optimizer, theoretically enabling the generation of complex few-shot examples.

B. Optimization Algorithm

We employed the DSPy framework with the BootstrapFewShot teleprompter. A critical challenge with small, quantized models is "hallucination," where the model generates plausible but incorrect reasoning steps. To mitigate this, we implemented a custom **Cost-Aware Metric** (Eq. 1). Unlike standard accuracy metrics, this metric penalizes answer verbosity:

$$\mathcal{L} = \mathcal{L}_{accuracy} \cdot e^{-\lambda \cdot length} \quad (1)$$

where $\lambda = 0.1$. This aggressive decay function forces the optimizer to prefer concise, direct answers, filtering out "rambling" reasoning traces that often indicate model confusion.

III. EXPERIMENTAL RESULTS

The evaluation was performed on the SQuAD (Stanford Question Answering Dataset) validation split ($N = 20$).

TABLE I
PERFORMANCE COMPARISON: ZERO-SHOT VS. OPTIMIZED

Stage	Accuracy	Qualitative Analysis
Baseline (Zero-Shot)	85.0%	Concise, extracted directly from text.
Optimized (Few-Shot)	80.0%	Minor hallucinations in reasoning.

A. Analysis of Regression

We observed a 5% regression (85% → 80)

The Llama-3.2-3B model is heavily instruction-tuned. When the optimizer injected few-shot examples (demonstrations) into the prompt, it consumed valuable tokens in the limited context window. For a 4-bit quantized model, the "signal-to-noise" ratio decreased; the model struggled to distinguish between the few-shot examples and the actual query. The zero-shot baseline, being cleaner and less "noisy," allowed the model to leverage its pre-trained instruction following capabilities more effectively.

IV. CONCLUSION

We successfully reproduced the software architecture of Automated Prompt Optimization on ultra-low-resource hardware, proving technical feasibility. However, our findings suggest a "complexity threshold": for models under 3 billion parameters, **Zero-Shot prompting is superior to automated Few-Shot optimization.** The overhead of processing in-context examples outweighs their instructional value at this scale. Future research for low-resource environments should prioritize "Instruction Tuning" optimizers (like MIPROv2) which refine the system prompt itself, rather than cluttering the context window with examples.

REFERENCES

- [1] Promptomatix Paper (ArXiv): <https://arxiv.org/abs/2507.14241>
- [2] Reproduction Logs (Notebook): https://colab.research.google.com/drive/Closest_Reproduced_Result.ipynb