

Instruction Fine-Tuning of LLaMA-2-7B with LoRA on Dolly-15K

Nauman (AUTHOR)¹

¹Department in Faculty of Computing, National University of Singapore, Singapore

nauman843@hotmail.com

Abstract – We fine-tune meta-llama/Llama-2-7b-hf using Parameter-Efficient Fine-Tuning (LoRA) on the databricks/databricks-dolly-15k dataset to adapt the base model for instruction following. We report training/validation loss curves and evaluate with AlpacaEval 2 (GPT-4 as judge over the AlpacaEval dataset) and MT-Bench (FastChat). Due to budget constraints, MT-Bench judgments are computed on a subset; full model answers for all 80 conversations are generated and stored for reproducibility. We document a known limitation: the single-turn Dolly format mismatches MT-Bench’s multi-turn conversational format, and the model may append ”Response:” artifacts. We provide cleaned, GitHub-renderable notebooks and scripts to reproduce fine-tuning, inference, and evaluation.

1 Background and Motivation

Large language models (LLMs) such as LLaMA-2-7B are powerful but require alignment to follow user instructions reliably. Instruction fine-tuning adapts a pretrained model from generic next-token prediction to helpful, safe, and coherent task completion. Parameter-Efficient Fine-Tuning (PEFT) with LoRA enables this alignment on modest hardware by updating a small set of low-rank adapter weights while keeping base weights frozen.

We target Dolly-15K (instruction, context, response) to train an assistant-style model and evaluate with two community benchmarks expected by the assignment: AlpacaEval 2 (instruction following, GPT-4 judged) and MT-Bench (multi-turn dialogue quality, GPT-4 judged). The goal is to demonstrate loss convergence, measurable gains over the base model, and to document practical limitations and reproducibility.

2 Relevance and related works

LoRA/PEFT has become a standard approach for efficient adaptation of LLMs, enabling fine-tuning on consumer-grade GPUs. Dolly-15K provides a permissive instruction-following corpus. AlpacaEval 2 and MT-Bench (via FastChat) are widely used evaluation protocols; both rely on LLM-as-a-judge (e.g., GPT-4) for scalable scoring. This work follows those best practices, while explicitly noting dataset–evaluation format mismatches (single-turn vs. multi-turn) as an important confounder for MT-Bench.

3 Methods and Approach

Dataset. We use databricks/databricks-dolly-15k. Each record is formatted as:

```
### Instruction:
{instruction}
```

```
### Context:
{context}
```

```
### Response:
{response}
```

We remove empty responses and split into train/validation/test (80/10/10).

Model

Training. Base: meta-llama/Llama-2-7b-hf. We apply LoRA with PEFT, enabling 4-bit quantization if needed for memory. Key settings include learning rate, epochs, global batch size, max length, gradient checkpointing, and saving LoRA adapters. Training/validation loss is tracked per epoch. Hardware details are recorded (GPU, VRAM, runtime).

Inference. For testing/inference, LoRA adapters are merged into the base model to create a standalone checkpoint for evaluation and answer generation.

Evaluation.

- *AlpacaEval 2:* It is a dataset of instruction-following prompts. We generate model responses and compute win rate using GPT-4 as an automated judge against reference responses.
- *MT-Bench (FastChat):* We generate model answers for 80 multi-turn conversations for both baseline and fine-tuned models. GPT-4 single-answer grading is run locally on our JSONL outputs (subset if budget-limited), producing overall and per-category scores. We document that the training prompt style may induce formatting artifacts (e.g., appending “ Response:”), which negatively impacts MT-Bench.

Reproducibility. All notebooks include explanatory markdown; Colab-specific metadata (e.g., `metadata.widgets`) is stripped to ensure GitHub rendering while preserving cell outputs. A FastChat submodule is added to run MT-Bench lo-

cally; answer files reside under the expected `model_answer/` directory and judgments are saved under `model_judgment/`.

4 Results and Evaluation

Artifacts and Reuse. The project delivers reproducible notebooks, scripts, and a cleaned repository ready for GitHub rendering, along with documented steps to re-run MT-Bench judging locally.

4.1 Training Dynamics

The LoRA fine-tuning process showed successful convergence over 751 training steps (approximately 3.2 hours). Training was conducted with the following key parameters:

- **Hardware:** NVIDIA A100 GPU (40GB VRAM)
- **Training Duration:** 194.59 minutes (3.2 hours)
- **Total Steps:** 751 steps
- **Final Training Loss:** 1.32 (converged from initial 1.45)
- **Learning Rate:** $2e-4$ (constant scheduler)
- **Batch Size:** 4 per device with 4 gradient accumulation steps
- **LoRA Parameters:** 39.9M trainable parameters (0.59% of total model)

Loss Convergence Analysis. The training loss decreased from approximately 1.45 at the beginning to 1.03 at the final step, demonstrating successful convergence. The loss curve shows the expected pattern of initial rapid decrease followed by gradual stabilization, indicating effective learning without overfitting.

4.2 AlpacaEval 2 Results

We evaluated the fine-tuned model against the baseline using AlpacaEval 2 with GPT-4 Turbo as judge. The results demonstrate significant improvements from LoRA fine-tuning:

Table 1: AlpacaEval 2 Results (GPT-4 Turbo judged)

Metric	Value
Win Rate	76.74%
Length-Controlled Win Rate	87.49%
Standard Error	13.02%
Average Response Length	167 tokens

The fine-tuned model (Llama2-7B-Dolly-QLoRA) achieved a **76.74

4.3 MT-Bench Results

We evaluated both the baseline and fine-tuned models on MT-Bench using GPT-4 as judge. The results show clear improvements from fine-tuning:

Table 2: MT-Bench Results (GPT-4 judged)

Model	Turn 1	Turn 2	Average
llama-2-7b-dolly-qlora	1.82	1.32	1.60
llama-2-7b-hf-baseline	1.48	1.24	1.38

The fine-tuned model (llama-2-7b-dolly-qlora) achieved a 16% improvement over the baseline (1.60 vs 1.38 average score). Both models show the expected pattern where second-turn scores are lower than first-turn scores, reflecting the increased difficulty of multi-turn conversations. The improvement is consistent across both turns, demonstrating that the fine-tuning successfully enhanced the model’s conversational capabilities.

5 Conclusion

This work successfully demonstrates the effectiveness of Parameter-Efficient Fine-Tuning (LoRA) for adapting LLaMA-2-7B to instruction-following tasks. The transformation from a generic language model to a helpful conversational assistant is evident in both quantitative metrics and qualitative behavior.

Key Achievements. Our LoRA fine-tuning approach achieved significant improvements across multiple evaluation dimensions. The fine-tuned model achieved a 76.74

Behavioral Transformation. The most striking evidence of successful fine-tuning lies in the model’s behavioral change. Before fine-tuning, when asked “What is the capital of France?”, the base model would continue with generic token generation, producing responses like “What is France?” or “What is Paris?” - essentially continuing the pattern rather than answering the question. After LoRA fine-tuning on Dolly-15K, the same prompt elicits a direct, helpful response: “Paris.”

Technical Efficiency. The LoRA approach proved highly efficient, training only 0.59

Reproducibility. All code, notebooks, and evaluation results are made available for reproducibility. The project demonstrates a complete pipeline from data preprocessing through training to evaluation, providing a template for similar instruction-following adaptations.

6 Data Policy

We open-source the code, notebooks, and configuration to reproduce training, inference, and evaluation. Notebooks are cleaned to render on GitHub while preserving outputs. MT-Bench model answers and (partial) GPT-4 judgments are saved in the repository structure expected by FastChat for local inspection. Dataset usage follows respective licenses.