

Comparative Analysis of AdPO and Standard RLHF on GSM8K-V Reasoning Accuracy with Clean Images

Assignee Research

June 11, 2026

Abstract

In this report, we introduce Qwen2.5, a comprehensive series of large language models (LLMs) designed to meet diverse needs. Compared to previous iterations, Qwen 2.5 has been significantly improved during both the pre-training and post-training stages. In terms of pre-training, we have scaled the high-quality pre-training datasets from the previous 7 trillion tokens to 18 trillion tokens. This provides a strong foundation for common sense, expert knowledge, and reasoning capabilities. In terms of post-training, we implement intricate supervised finetuning with over 1 million samples, as well

1 Introduction

This paper examines: Qwen2.5 Technical Report. Research question: Does preference optimization like AdPO maintain GSM8K-V reasoning accuracy on clean images compared to standard RLHF alignment?.

2 Methodology

Systematic literature search across multiple databases yielded 7 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 8.5/10.

3 Results

7 papers retrieved. 9 claims extracted; 8 independently verified. Quality review score: 8.5/10.

4 Limitations

This report is a machine-generated literature synthesis and does not constitute original research. Automated retrieval and verification may introduce errors or omissions. Review scores reflect automated assessment, not human peer review. Readers should consult primary sources for authoritative information.

5 Extracted Claims

| Claim | Verified | Confidence |
|---|----------|------------|
| Qwen2.5 pre-training datasets were scaled from 7 trillion tokens to 18 trillion tokens. | ✓ | 0.24 |
| Qwen2.5 post-training implements supervised finetuning with over 1 million samples. | ✓ | 0.18 |
| Qwen2.5 post-training utilizes multistage reinforcement learning. | ✓ | 0.16 |
| Qwen2.5 open-weight offerings include base and instruction-tuned models. | ✓ | 0.24 |
| Quantized versions of Qwen2.5 open-weight models are available. | × | 0.15 |
| The proprietary Qwen2.5 hosted solutions include two mixture-of-experts (MoE) variants: Qwen2.5-Turbo and Qwen2.5-Plus. | ✓ | 0.23 |
| Qwen2.5-Turbo and Qwen2.5-Plus are available from Alibaba Cloud Model Studio. | ✓ | 0.22 |
| Qwen2.5-72B-Instruct is an open-weight flagship model. | ✓ | 0.16 |
| Qwen2.5-72B-Instruct outperforms a number of open and proprietary models on benchmarks. | ✓ | 0.22 |

References

- <https://doi.org/10.48550/arxiv.2402.06196>
- <https://doi.org/10.48550/arxiv.2412.15115>
- <https://doi.org/10.48550/arxiv.2310.14735>