

Long-context reasoning impact on code generation performance in synthetic scenarios

Assignee Research

June 12, 2026

Abstract

Adapting large language models (LLMs) to long-context tasks requires post-training methods that remain accurate and coherent over thousands of tokens. Existing approaches are limited in several ways: 1) off-policy methods such as supervised fine-tuning (SFT) and knowledge distillation (KD) suffer from exposure bias and limited recovery from model-generated errors over long horizons; 2) on-policy reinforcement learning methods such as Group Relative Policy Optimization (GRPO) better align training with model-generated states, but are unstable and sample-inefficient due to sparse rewards; 3) on-

1 Introduction

This paper examines: Combining On-Policy Optimization and Distillation for Long-Context Reasoning in Large Language Models. Research question: Does the long-context reasoning improvement observed in LongReason benchmarks translate to improved performance on code generation tasks in synthetic long-context scenarios, and how does it compare to short-context code generation accuracy?.

2 Methodology

Systematic literature search across multiple databases yielded 8 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 8.0/10.

3 Results

8 papers retrieved. 14 claims extracted; 13 independently verified. Quality review score: 8.0/10.

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
Off-policy methods such as supervised fine-tuning (SFT) and knowledge distillation (KD) suffer from exposure bias over l	✓	0.30
Off-policy methods such as SFT and KD have limited recovery from model-generated errors over long horizons.	✓	0.24
On-policy reinforcement learning methods such as Group Relative Policy Optimization (GRPO) better align training with mo	✓	0.35
On-policy reinforcement learning methods such as GRPO are unstable due to sparse rewards.	✓	0.17
On-policy reinforcement learning methods such as GRPO are sample-inefficient due to sparse rewards.	✓	0.21
On-policy distillation (OPD) provides dense token-level guidance.	✓	0.25
On-policy distillation (OPD) does not directly optimize arbitrary reward signals.	✓	0.24
The paper proposes Distilled Group Relative Policy Optimization (dGRPO), a method that augments GRPO with dense guidance	✓	0.32
The paper introduces LongBlocks, a synthetic long-context dataset.	✓	0.18
The LongBlocks dataset spans multi-hop reasoning, contextual grounding, and long-form generation.	✓	0.21
Experiments in the paper compare off-policy training, sparse-reward GRPO, and the combined dGRPO approach.	✓	0.19
Combining outcome-based policy optimization with knowledge distillation in a single objective provides a more stable pat	✓	0.33
Combining outcome-based policy optimization with knowledge distillation in a single objective provides a more effective	✓	0.33
The proposed dGRPO method preserves short-context capabilities.	×	0.11

References

- <https://openalex.org/W7163597136>
- <https://doi.org/10.1007/s10489-026-07230-0>
- <https://openalex.org/W7161205011>