

# Test-Time Compute Scaling and Its Impact on Language Model Reasoning Performance

Assignee Research

June 6, 2026

## Abstract

This report synthesises findings from 13 peer-reviewed papers addressing the following research question: How does test-time compute scaling improve language model performance on reasoning benchmarks v17. 13 claims were extracted from source literature; 1 was independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.5/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: ARS: Adaptive Reasoning Suppression for Efficient Large Reasoning Language Models. Research question: How does test-time compute scaling improve language model performance on reasoning benchmarks v17.

## 2 Methodology

Systematic literature search across multiple databases yielded 13 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 4.5/10.

## 3 Results

13 papers retrieved. 13 claims extracted; 1 independently verified. Quality review score: 4.5/10.

## 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
The study evaluates Qwen2.5-Math-1.5B-Instruct, Qwen2.5-Math-7B-Instruct, and DeepSeek-R1-Distill-Qwen-7B models.	×	0.04
The ARS algorithm utilizes difficulty thresholds d1 and d2 to determine the scheduling mode.	×	0.03
In FAST mode, the ARS policy is configured with 2 drafts and 10 tokens per draft.	×	0.01
In MOD mode, the ARS policy uses a budget of 64 tokens.	×	0.02
The generation process in the ARS algorithm terminates if the text length reaches 1200 tokens.	×	0.02
Confidence scores in ARS are computed using entropy confidence of a tentative answer.	×	0.04
Token suppression occurs if the next token is in the trigger set and the suppression probability exceeds a random value.	×	0.07
The trigger set T for reflection behaviors includes keywords such as 'Wait', 'But', and 'Alternatively'.	×	0.01
The objective of the ARS framework is to minimize expected output length $E[T]$ while keeping accuracy degradation below a	×	0.02
ARS operates through three core components: Multi-checkpoint certainty estimation, Progressive threshold adaptation, and	✓	0.15
Table 1 presents performance comparison results on the GSM8K dataset.	×	0.04
Table 2 presents performance comparison results on the MATH500 dataset.	×	0.04
The text claims that ARS consistently achieves superior length reduction while maintaining competitive accuracy across a	×	0.09

## References

- <http://arxiv.org/abs/2601.01982v1>
- <http://arxiv.org/abs/2510.00071v2>
- <http://arxiv.org/abs/2504.13171v1>