

# Test-Time Compute Scaling Enhances Reasoning Benchmarks in Sub-10B Language Models

Assignee Research

June 6, 2026

## Abstract

This report synthesises findings from 16 peer-reviewed papers addressing the following research question: How does test-time compute scaling improve language model performance on reasoning benchmarks v11. 19 claims were extracted from source literature; 7 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 5.8/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: m1: Unleash the Potential of Test-Time Scaling for Medical Reasoning with Large Language Models. Research question: How does test-time compute scaling improve language model performance on reasoning benchmarks v11.

## 2 Methodology

Systematic literature search across multiple databases yielded 16 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 5.8/10.

## 3 Results

16 papers retrieved. 19 claims extracted; 7 independently verified. Quality review score: 5.8/10.

## 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
Test-time scaling (by increasing the 'thinking' token budget) consistently enhances medical reasoning.	✓	0.28
Lightweight fine-tuned models under 10B parameters establish new state-of-the-art performance with test-time scaling.	✓	0.31
A 32B model achieves results comparable to previous 70B-scale medical LLMs.	✓	0.16
There is an optimal reasoning token budget of approximately 4K, beyond which performance may degrade due to overthinking	✓	0.24
Budget forcing does not necessarily improve overall medical QA performance and can introduce errors into previously corr	✓	0.25
Insufficient medical knowledge is a key bottleneck that prevents further performance gains through test-time scaling.	✓	0.32
Increasing data scale, improving data quality, and expanding model capacity enhance medical knowledge grounding.	✓	0.24
Enriched medical knowledge is essential for fully realizing the benefits of test-time scaling.	×	0.13
The code, models, and data are publicly available at <a href="https://github.com/UCSC-VLAA/m1">https://github.com/UCSC-VLAA/m1</a> .	×	0.09
The research does not require IRB approval.	×	0.06
The evaluation includes nine medical QA benchmarks, grouped into In-Distribution and Out-of-Distribution tests.	×	0.04
Accuracy is measured for all datasets.	×	0.03
In-Distribution tests include MedMCQA, MedQA-USMLE, and PubMedQA.	×	0.02
Out-of-Distribution tests include MMLU-Pro, GPQA, Lancet, NEJM, MedBullets (4 Options and 5 Options), and MedXpertQA.	×	0.03
The models are compared against a variety of general and specialized medical LLM benchmarks.	×	0.04
The performance of models consistently improves with increased inference.	×	0.06
The table shows the performance of various models on different medical QA benchmarks.	×	0.06
The table shows the performance of various models on different medical QA benchmarks with data filtering.	×	0.06
The table shows the performance of various models on different medical QA benchmarks with data filtering.	×	0.06

## References

- <http://arxiv.org/abs/2505.03786v1>
- <http://arxiv.org/abs/2502.05171v2>
- <http://arxiv.org/abs/2504.00869v2>