

# GLM-4-9B Benchmark Performance Across Reasoning Mathematics Coding and Language Tasks

Assignee Research

June 6, 2026

## Abstract

This report synthesises findings from 10 peer-reviewed papers addressing the following research question: What are the benchmark performance scores of GLM-4-9B on reasoning mathematics coding and language understanding tasks. 13 claims were extracted from source literature; 4 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 5.3/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: Qwen2 Technical Report. Research question: What are the benchmark performance scores of GLM-4-9B on reasoning mathematics coding and language understanding tasks.

## 2 Methodology

Systematic literature search across multiple databases yielded 10 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 5.3/10.

## 3 Results

10 papers retrieved. 13 claims extracted; 4 independently verified. Quality review score: 5.3/10.

## 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 Qwen2 series includes foundational and instruction-tuned language models with parameter counts ranging from 0.5 bill	✓	0.18
The Qwen2 series features both dense models and a Mixture-of-Experts model.	✓	0.17
Qwen2-72B achieves a score of 84.2 on the MMLU benchmark.	×	0.09
Qwen2-72B achieves a score of 37.9 on the GPQA benchmark.	×	0.09
Qwen2-72B achieves a score of 64.6 on the HumanEval benchmark.	×	0.09
Qwen2-72B achieves a score of 89.5 on the GSM8K benchmark.	×	0.09
Qwen2-72B achieves a score of 82.4 on the BBH benchmark.	×	0.09
Qwen2-72B-Instruct achieves a score of 9.1 on the MT-Bench benchmark.	×	0.11
Qwen2-72B-Instruct achieves a score of 48.1 on the Arena-Hard benchmark.	×	0.13
Qwen2-72B-Instruct achieves a score of 35.7 on the LiveCodeBench benchmark.	×	0.11
Qwen2 is proficient in approximately 30 languages.	×	0.14
Qwen2 model weights are available on Hugging Face and ModelScope.	✓	0.18
Supplementary materials including example code for Qwen2 are available on GitHub.	✓	0.17

## References

- <https://doi.org/10.48550/arxiv.2406.12793>

- <https://doi.org/10.18653/v1/2025.findings-emnlp.20>
- <https://doi.org/10.48550/arxiv.2407.10671>