

SOVEREIGN: Claude-3 evaluation benchmark MMLU HumanEval GSM8K MATH coding performance scores Anthropic

SOVEREIGN Research Kernel

Autonomous draft — Owner review required before publication

May 29, 2026

Abstract

We introduce phi-3-mini, a 3.8 billion parameter language model trained on 3.3 trillion tokens, whose overall performance, as measured by both academic benchmarks and internal testing, rivals that of models such as Mixtral 8x7B and GPT-3.5 (e.g., phi-3-mini achieves 69% on MMLU and 8.38 on MT-bench), despite being small enough to be deployed on a phone. Our training dataset is a scaled-up version of the one used for phi-2, composed of heavily filtered publicly available web data and synthetic data. The model is also further aligned for robustness, safety, and chat format. We also provide param

1 Introduction

Analysis of: Phi-3 Technical Report: A Highly Capable Language Model Locally on Your Phone. Research goal: Claude-3 evaluation benchmark MMLU HumanEval GSM8K MATH coding performance scores Anthropic.

2 Methodology

Multi-query arXiv search (4 parallel queries, Relevance-sorted). TF-IDF cosine semantic verification (bigrams, threshold=0.15). NIM nv-embedqa-e5-v5 (dim=1024) for semantic indexing. Tribunal v2: 3-role parallel review (SKEPTIC/VALIDATOR/SYNTHESIZER) with revision round if score < 6.5.

3 Results

9 papers retrieved. 14 claims extracted, 14 verified. Tribunal: 9.2/10 \rightarrow APPROVE (revision_round=0). Policy: AUTO_APPROVE.

4 Uncertainties

NIM free tier latency varies. TF-IDF verification is a weak signal. arXiv Relevance ranking is query-dependent. Tribunal consensus is LLM-based and prompt-sensitive.

5 Extracted Claims

Claim	Verified	Confidence
phi-3-mini is a 3.8 billion parameter language model trained on 3.3 trillion tokens.	✓	0.29
phi-3-mini achieves 69% on MMLU and 8.38 on MT-bench.	✓	0.26
phi-3-mini’s performance rivals that of models such as Mixtral 8x7B and GPT-3.5.	✓	0.21
phi-3-mini is small enough to be deployed on a phone.	✓	0.19
The training dataset for phi-3-mini is a scaled-up version of the one used for phi-2, composed of heavily filtered publi	✓	0.31
phi-3-mini is further aligned for robustness, safety, and chat format.	✓	0.20
phi-3-small is a 7B model trained for 4.8T tokens, achieving 75% on MMLU and 8.7 on MT-bench.	✓	0.21
phi-3-medium is a 14B model trained for 4.8T tokens, achieving 78% on MMLU and 8.9 on MT-bench.	✓	0.22
phi-3.5-mini is part of the phi-3.5 series introduced to enhance multilingual, multimodal, and long-context capabilities	✓	0.23
phi-3.5-MoE is a 16 x 3.8B MoE model with 6.6 billion active parameters.	✓	0.26
phi-3.5-MoE achieves superior performance in language reasoning, math, and code tasks compared to other open-source mode	✓	0.34
phi-3.5-MoE performs on par with Gemini-1.5-Flash and GPT-4o-mini.	✓	0.20
phi-3.5-Vision is a 4.2 billion parameter model derived from phi-3.5-mini.	✓	0.29
phi-3.5-Vision excels in reasoning tasks and is adept at handling both single-image and text prompts, as well as multi-i	✓	0.33

References

- <https://doi.org/10.48550/arxiv.2406.11931>
- <https://doi.org/10.48550/arxiv.2403.05530>
- <https://doi.org/10.48550/arxiv.2404.14219>