

SOVEREIGN: What is the performance gap in code generation tasks (HumanEval) between multimodal models and text-only LLMs

SOVEREIGN Research Kernel

Autonomous draft — Owner review required before publication

May 29, 2026

Abstract

We release Code Llama, a family of large language models for code based on Llama 2 providing state-of-the-art performance among open models, infilling capabilities, support for large input contexts, and zero-shot instruction following ability for programming tasks. We provide multiple flavors to cover a wide range of applications: foundation models (Code Llama), Python specializations (Code Llama - Python), and instruction-following models (Code Llama - Instruct) with 7B, 13B, 34B and 70B parameters each. All models are trained on sequences of 16k tokens and show improvements on inputs with up

1 Introduction

Analysis of: Code Llama: Open Foundation Models for Code. Research goal: What is the performance gap in code generation tasks (HumanEval) between multimodal models and text-only LLMs when provided with additional visual context or pseudocode diagrams?.

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

14 papers retrieved. 16 claims extracted, 14 verified. Tribunal: 7.7/10 \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
Code Llama is a family of large language models for code based on Llama 2.	✓	0.28
Code Llama provides state-of-the-art performance among open models.	✓	0.25
Code Llama has infilling capabilities.	✓	0.16
Code Llama supports large input contexts.	×	0.15
Code Llama has zero-shot instruction following ability for programming tasks.	✓	0.28
Code Llama is available in multiple flavors: foundation models (Code Llama), Python specializations (Code Llama - Python	✓	0.41
Code Llama models have 7B, 13B, 34B, and 70B parameters each.	✓	0.27
All Code Llama models are trained on sequences of 16k tokens.	✓	0.24
Code Llama models show improvements on inputs with up to 100k tokens.	✓	0.22
7B, 13B, and 70B Code Llama and Code Llama - Instruct variants support infilling based on surrounding content.	✓	0.41
Code Llama reaches state-of-the-art performance among open models on several code benchmarks.	✓	0.36
Code Llama scores up to 67% on HumanEval.	✓	0.17
Code Llama scores up to 65% on MBPP.	×	0.14
Code Llama - Python 7B outperforms Llama 2 70B on HumanEval and MBPP.	✓	0.36
All Code Llama models outperform every other publicly available model on MultiPL-E.	✓	0.24
Code Llama is released under a permissive license that allows for both research and commercial use.	✓	0.21

References

- <https://doi.org/10.1039/d4sc03921a>
- <https://doi.org/10.48550/arxiv.2303.12712>
- <https://doi.org/10.48550/arxiv.2308.12950>