

# Language Models in Formal Theorem Proving and Mathematical Verification Tasks

Assignee Research

June 6, 2026

## Abstract

This report synthesises findings from 20 peer-reviewed papers addressing the following research question: How do language models perform on formal theorem proving and mathematical verification tasks v17. 15 claims were extracted from source literature; 9 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 7.2/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: Automated Identification of Incidentalomas Requiring Follow-Up: A Multi-Anatomy Evaluation of LLM-Based and Supervised Approaches. Research question: How do language models perform on formal theorem proving and mathematical verification tasks v17.

## 2 Methodology

Systematic literature search across multiple databases yielded 20 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 7.2/10.

## 3 Results

20 papers retrieved. 15 claims extracted; 9 independently verified. Quality review score: 7.2/10.

## 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 anatomy-informed GPT-OSS-20b model achieved the highest performance, yielding an incidentaloma-positive macro-F1 of	✓	0.32
This surpassed all supervised baselines (maximum macro-F1: 0.70) and closely matched the inter-annotator agreement of 0.	✓	0.30
Explicit anatomical grounding yielded statistically significant performance gains across GPT-based models ( $p < 0.05$ ).	✓	0.27
A majority-vote ensemble of the top systems further improved the macro-F1 to 0.90.	✓	0.22
Error analysis revealed that anatomy-aware LLMs demonstrated superior contextual reasoning in distinguishing actionable	✓	0.31
Generative LLMs, when enhanced with structured lesion tagging and anatomical context, significantly outperform tradition	✓	0.26
The best overall performance was achieved by the GPT-OSS-20b (With Anatomy) model, which obtained the highest F1-scores	✓	0.20
GPT-4o (With Anatomy) was the second-best performer, with F1 scores of 0.82 and 0.71 for the incidentaloma-positive clas	×	0.10
These two anatomy-informed methods consistently outperformed all other systems, indicating that explicit anatomy context	×	0.08
Among supervised encoders, BioClinicalModernBERT (w/o CS) and ModernBERT (CS) both reached an incidentaloma macro-F1 of	×	0.13
Cost-sensitive (CS) learning produced only modest gains and mainly increased recall for minority classes.	×	0.02
Comparing model families, LLMs showed clear gains over supervised encoders.	×	0.08
Only Llama 3.1-8B was fine-tuned using LoRA; all other LLMs, including GPT-4o and GPT-OSS-20B, were evaluated in a promp	✓	0.15
GPT-4o (Base) already matched the strongest non-LLM baselines, and adding anatomical context further improved performanc	×	0.14
The consistent benefit of anatomy-aware prompting across both Llama and GPT-based architectures highlights the value of 4	✓	0.15

## References

- <https://www.semanticscholar.org/paper/f52e524e719054f7c24d81dcf153edca914145ab>
- <https://arxiv.org/abs/2512.05537>
- <https://www.semanticscholar.org/paper/22789c1c946ab0ced8702813ca4b7b64a3819554>