

SOVEREIGN: What is the impact of test-time compute allocation on the inference efficiency and task completion accuracy of

SOVEREIGN Research Kernel

Autonomous draft — Owner review required before publication

May 29, 2026

Abstract

Recent advances in test-time scaling of large language models (LLMs), exemplified by DeepSeek-R1 and OpenAI’s o1, show that extending the chain of thought during inference can significantly improve general reasoning performance. However, the impact of this paradigm on legal reasoning remains insufficiently explored. To address this gap, we present the first systematic evaluation of 12 LLMs, including both reasoning-focused and general-purpose models, across 17 Chinese and English legal tasks spanning statutory and case-law traditions. In addition, we curate a bilingual chain-of-thought dataset

1 Introduction

Analysis of: Evaluating Test-Time Scaling LLMs for Legal Reasoning: OpenAI o1, DeepSeek-R1, and Beyond. Research goal: What is the impact of test-time compute allocation on the inference efficiency and task completion accuracy of DeepSeek-R1 versus o1-preview models in multilingual legal document classification tasks.

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

7 papers retrieved. 5 claims extracted, 2 verified. Tribunal: 5.0/10 \rightarrow REVERSE (revision_round=1). Policy: ESCALATE_TO_OWNER.

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
DeepSeek-R1 demonstrates superior performance in Chinese legal reasoning tasks.	✓	0.19
Legal-R1 outperforms baseline models on the majority of legal tasks.	×	0.09
Legal-R1-14B was developed through progressive supervised fine-tuning and achieves enhanced performance on legal tasks.	×	0.07
The dataset and model are available at https://github.com/YinghaoHu/Legal-R1-14B	×	0.06
Legal reasoning performance is lacking in current LLMs, with main obstacles being outdated legal knowledge and hallucina	✓	0.17

References

- <http://arxiv.org/abs/2509.22472v1>
- <http://arxiv.org/abs/2310.05276v1>
- <http://arxiv.org/abs/2503.16040v2>