

# SOVEREIGN: To what extent does ReSSFormer’s sparse attention pattern improve long-context generalization on NLVR2 over de

SOVEREIGN Research Kernel

Autonomous draft — Owner review required before publication

May 28, 2026

## Abstract

Foundation models, now powering most of the exciting applications in deep learning, are almost universally based on the Transformer architecture and its core attention module. Many subquadratic-time architectures such as linear attention, gated convolution and recurrent models, and structured state space models (SSMs) have been developed to address Transformers’ computational inefficiency on long sequences, but they have not performed as well as attention on important modalities such as language. We identify that a key weakness of such models is their inability to perform content-based reasoni

## 1 Introduction

Analysis of: Mamba: Linear-Time Sequence Modeling with Selective State Spaces. Research goal: To what extent does ReSSFormer’s sparse attention pattern improve long-context generalization on NLVR2 over dense transformers when evaluated on sequences exceeding 4K tokens, measured by exact match scores and memory usage?.

## 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

5 papers retrieved. 3 claims extracted, 3 verified. Tribunal: 7.8/10 → 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
Mamba achieves 5× higher throughput than Transformers	✓	0.15
Mamba achieves state-of-the-art performance across several modalities such as language, audio, and genomics	✓	0.24
Mamba enjoys fast inference and linear scaling in sequence length	✓	0.25

### References

- <https://doi.org/10.1186/s40537-021-00444-8>
- <https://doi.org/10.48550/arxiv.2312.00752>
- <https://doi.org/10.48550/arxiv.2302.13971>