

SOVEREIGN: How does the choice of attention mechanism (e.g., sparse vs. dense) in vision transformers affect mean Interse

SOVEREIGN Research Kernel

Autonomous draft — Owner review required before publication

May 29, 2026

Abstract

Since the introduction of Vision Transformers, the landscape of many computer vision tasks (e.g., semantic segmentation), which has been overwhelmingly dominated by CNNs, recently has significantly revolutionized. However, the computational cost and memory requirement renders these methods unsuitable on the mobile device. In this paper, we introduce a new method squeeze-enhanced Axial Transformer (SeaFormer) for mobile visual recognition. Specifically, we design a generic attention block characterized by the formulation of squeeze Axial and detail enhancement. It can be further used to create

1 Introduction

Analysis of: SeaFormer++: Squeeze-enhanced Axial Transformer for Mobile Visual Recognition. Research goal: How does the choice of attention mechanism (e.g., sparse vs. dense) in vision transformers affect mean Intersection over Union (mIoU) on driving scene segmentation benchmarks (Cityscapes, BDD100K) under real-time latency constraints?.

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

8 papers retrieved. 11 claims extracted, 10 verified. Tribunal: 8.2/10
\$\\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
Vision Transformers have significantly revolutionized computer vision tasks previously dominated by CNNs, such as semant	✓	0.19
The computational cost and memory requirements of existing Vision Transformer methods render them unsuitable for mobile	✓	0.16
SeaFormer is a squeeze-enhanced Axial Transformer designed for mobile visual recognition.	✓	0.31
The SeaFormer attention block is characterized by the formulation of squeeze Axial and detail enhancement.	✓	0.23
SeaFormer can be used to create a family of backbone architectures with superior cost-effectiveness.	✓	0.23
When coupled with a light segmentation head, SeaFormer achieves the best trade-off between segmentation accuracy and lat	✓	0.28
SeaFormer was evaluated on the ADE20K, Cityscapes, Pascal Context, and COCO-Stuff datasets.	✓	0.19
SeaFormer outperforms both mobile-friendly rivals and Transformer-based counterparts in terms of performance and latency	×	0.15
A feature upsampling-based multi-resolution distillation technique was incorporated into the SeaFormer framework to redu	✓	0.20
The SeaFormer architecture was applied to image classification and object detection problems beyond semantic segmentatio	✓	0.20
The code and models for SeaFormer are publicly available at https://github.com/fudan-zvg/SeaF .	✓	0.22

References

- <https://doi.org/10.1109/tits.2022.3207665>
- <https://doi.org/10.48550/arxiv.2301.13156>
- <https://doi.org/10.3390/electronics12122730>