

SOVEREIGN: How does the proposed CAT method compare to fixed tokenization approaches in terms of end-to-end inference time

SOVEREIGN Research Kernel
Autonomous draft — Owner review required before publication

May 27, 2026

Abstract

Transformer-based video diffusion models rely on 3D attention over spatial and temporal tokens, which incurs quadratic time and memory complexity and makes end-to-end training for ultra-high-resolution videos prohibitively expensive. To overcome this bottleneck, we propose a pure image adaptation framework that upgrades a video Diffusion Transformer pre-trained at its native scale to synthesize higher-resolution videos. Unfortunately, naively fine-tuning with high-resolution images alone often introduces noticeable noise due to the image-video modality gap. To address this, we decouple the learning

1 Introduction

Analysis of: ViBe: Ultra-High-Resolution Video Synthesis Born from Pure Images. Research goal: How does the proposed CAT method compare to fixed tokenization approaches in terms of end-to-end inference time and token count scaling across different domain shifts in image complexity?.

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

1 papers retrieved. 11 claims extracted, 11 verified. Tribunal: 7.5/10 → 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
Transformer-based video diffusion models rely on 3D attention over spatial and temporal tokens, which incurs quadratic t	✓	0.34
End-to-end training for ultra-high-resolution videos is prohibitively expensive due to quadratic complexity of 3D attent	✓	0.28
ViBe proposes a pure image adaptation framework that upgrades a video Diffusion Transformer pre-trained at its native sc	✓	0.35
Naively fine-tuning with high-resolution images alone introduces noticeable noise due to the image-video modality gap.	✓	0.35
ViBe decouples the learning objective to separately handle modality alignment and spatial extrapolation.	✓	0.23
Relay LoRA is a two-stage adaptation strategy at the core of ViBe.	✓	0.17
In the first stage of Relay LoRA, the video diffusion model is adapted to the image domain using low-resolution images t	✓	0.36
In the second stage of Relay LoRA, the model is further adapted with high-resolution images to acquire spatial extrapola	✓	0.32
During inference, only the high-resolution adaptation is retained to preserve the video generation modality while enabli	✓	0.36
ViBe proposes a High-Frequency-Awareness-Training-Objective that explicitly encourages the model to recover high-frequen	✓	0.35
Extensive experiments demonstrate that ViBe produces ultra-high-resolution videos.	✓	0.24

References

- <https://www.semanticscholar.org/paper/efaa2bddd99ea37255a6ae31bff7578707266a7d>