

Paper Critique: VideoMamba: State Space Model for Efficient Video Understanding

Dan Peng
Department of Computer Science
Sep 17, 2024
danpeng@unc.edu

1. Research Problem

1.1. What research problem does the paper address?

This paper addressed the problem of video understanding that current SOTA models (CNNs and transformers) cannot solve both the large spatiotemporal redundancy within short video clips and the complex spatiotemporal dependencies among long contexts simultaneously.

1.2. What is the motivation of the research work?

The motivation came from the outstanding performance of state space models (SSMs) in NLP and vision tasks, which offer the significant reduction in memory usage and the high efficiency in long-term sequence modeling with linear complexity.

2. Technical Novelty

2.1. What are the key technical challenges identified by the authors?

Key technical challenges are: reducing Spatiotemporal Redundancy, modeling long-term dependency, and improving computational and memory efficiency.

2.2. How significant is the technical contribution of the paper? If you think that the paper is incremental, please provide references to the most similar work

This paper offers a purely SSM-based model designed for video understanding by: 1) significant improvements in computational efficiency (VideoMamba operates 6× faster than TimeSformer and requires 40× less GPU memory for long video sequences); 2) a novel self-distillation technique that enhances scalability without requiring large-scale dataset pretraining. This paper is incremental to the most similar work Vision Mamba on ImageNet and TimeSformer and ViViT on Kinetics-400 and SthSth V2.

2.3. Identify 1-5 main strengths of the proposed approach.

- Introduces a linear-complexity operator that handles long-term dependencies with much lower computational and memory costs compared to transformer-based models
- Scale to high-resolution and long video sequences without requiring extensive dataset pretraining, thanks to the self-distillation technique

2.4. Identify 1-5 main weaknesses of the proposed approach.

- Larger variants of VideoMamba tend to overfit on smaller datasets
- Not fully explored multi-modal capabilities like incorporating audio or the scalability of VideoMamba to extremely large datasets

3. Empirical Results

3.1. Identify 1-5 key experimental results, and explain what they signify.

- VideoMamba is 6× faster and requires 40× less GPU memory than TimeSformer for long-video sequences (64 frames). This experimental result highlights the model's computational efficiency and ability to process long-term video sequences.
- VideoMamba (25M) achieves a 42.6% top-1 retrieval accuracy on MSRVT and 43.1% on DiDeMo, outperforming existing models such as UMT and VideoCLIP with fewer parameters. It indicates VideoMamba's multi-modal task performance, specifically video-text retrieval.

3.2. Are there any weaknesses in the experimental section (i.e., unfair comparisons, missing ablations, etc)?

Yes, there are two weaknesses in the experimental design. First, VideoMamba is not evaluated on larger, more complex video datasets like Kinetics-600, which may not fully validate the argument about VideoMamba's scalability. Second, some comparisons between VideoMamba and other models (e.g., TimeSformer, ViViT) involve different pretraining datasets. For instance, VideoMamba is pre-trained on IN-1K, while some of the compared models are pre-trained on larger datasets like IN-21K or CLIP-400M. This would make the comparison unreliable and unfair.

4. Summary

The VideoMamba paper presents an innovative approach to video understanding using SSMs, achieving linear complexity and exceptional performance on both short-term and long-term video tasks. It also reduced memory and computational costs compared to traditional transformer-based models. I like this paper because it introduced a self-distillation technique which further enhances scalability without requiring large pretraining datasets.

5. QA Prompt for a Paper Discussion

5.1. Discussion Question

Can VideoMamba be adapted for real-time autonomous driving systems? If so, how? If no, why?

5.2. Your Answer

Yes, VideoMamba can be adapted for real-time driving systems by its high efficiency and ability to model long-term spatiotemporal information with low memory. We can integrate real-time multi-modal inputs into VideoMamba, in addition to video, like LiDAR, radar, and audio to handle complex driving environments. Also, the SSM's linear complexity ensures that high-resolution, long-duration video streamings from cameras can be processed in real-time dynamic scenarios.