

# Paper Critique – VideoMAE: Masked Autoencoders are data-efficient Learners for Self-Supervised Video Pre-Training

Alessandro Folloni  
Computer Science  
afolloni@unc.edu

## 1. Research Problem

### 1.1. Research Problem Addressed

The paper faces the need of pre-training in order to achieve optimal performance on small datasets.

### 1.2. Motivation of the Research

The motivation of the research is to address the challenge of training effective vision transformers (ViTs) for video recognition without relying on pre-trained models from large-scale image data. The authors recognize that while ViTs have shown significant progress in natural language processing and computer vision tasks, training ViTs for video data from scratch has proven to be challenging. Existing video datasets are smaller compared to image datasets, and video transformers tend to inherit biases from image-based models. The research aims to develop a self-supervised video pre-training method that can effectively train ViTs on video datasets without the need for pre-trained models, thus reducing the reliance on image-based biases.

## 2. Technical Novelty

### 2.1. Key Technical Challenges Identified

The key technical challenge identified in the research is the effective training of vision transformers for video data from scratch. This challenge arises due to several factors: limited video data, inherited biases, information leakage, slowness in video data.

### 2.2. Significance of the Technical Contribution

The technical contribution of the research is highly significant for the following reasons:

- **Introducing VideoMAE:** The research presents a novel self-supervised video pre-training method called Video Masked Autoencoder (VideoMAE). VideoMAE addresses the challenges of training video transformers from scratch by introducing a customized design that

includes tube masking with extremely high masking ratios and temporal downsampling.

- **Reduced reliance on pre-trained models:** VideoMAE allows for the training of vanilla ViT backbones on relatively small-scale video datasets without the need for pre-trained models from large-scale image data, reducing biases and increasing efficiency.
- **Improved video representation learning:** VideoMAE demonstrates that the proposed masking and reconstruction strategy provides an effective solution for self-supervised video pre-training. Models pre-trained with VideoMAE outperform those trained from scratch or pre-trained using contrastive learning methods.
- **Data-efficient learning:** The research shows that VideoMAE can be successfully trained with a relatively small number of videos, emphasizing data efficiency as an essential aspect of self-supervised video pre-training.

### 2.3. Main Strengths of the Proposed Approach

- **Effective masking strategy:** the use of tube masking with extremely high masking ratios and temporal tube masking helps address the challenge of information leakage and encourages the learning of spatiotemporal structures.
- **Reduced reliance on external data:** the approach enables the training of video transformers without relying on pre-trained models from large-scale image data, making it more suitable for video-specific tasks.
- **Data efficiency:** the research demonstrates that VideoMAE can be trained effectively with a relatively small number of videos, highlighting its data-efficient learning capabilities.

### 2.4. Main Weaknesses of the Proposed Approach

- **Computational complexity:** While the extremely high masking ratio is beneficial, it may still pose computational challenges in terms of processing only a small

percentage of unmasked tokens during pre-training. The research mentions this as a potential issue but does not provide extensive details on how this challenge is addressed.

### 3. Empirical Results

#### 3.1. Key Experimental Results and Their Significance

- VideoMAE outperforms other pre-training methods on video datasets, such as Kinetics-400 and Something-Something V2, achieving top-1 accuracy of 87.4% on Kinetics-400 and 75.4% on Something-Something V2 without using any external data.
- VideoMAE demonstrates data efficiency, achieving competitive results even with a small number of training videos (e.g., 3.5k videos on HMDB51).
- Transferability of VideoMAE representations is shown by fine-tuning on downstream tasks like action detection (e.g., achieving 26.7 mAP on AVA) and achieving better results than models pre-trained with other methods.
- VideoMAE can scale up effectively with more powerful backbones and larger input resolutions, achieving even better performance (e.g., 87.4

#### 3.2. Weaknesses in the Experimental Section

One potential weakness in the experimental section is the lack of a detailed analysis of the computational resources required for pre-training with VideoMAE. The paper mentions the efficiency of VideoMAE in terms of training time, but it would be beneficial to provide more specific information on GPU hours or other computational metrics to give a better understanding of the resource requirements.

### 4. Summary

The paper presents VideoMAE, a self-supervised pre-training method for video transformers. VideoMAE achieves state-of-the-art results on video datasets, demonstrates data efficiency, and shows strong transferability to downstream tasks. It introduces novel designs like an extremely high masking ratio and tube masking strategy to make video reconstruction more challenging. While the paper provides valuable contributions to the field, it could benefit from a more detailed analysis of computational resources used during pre-training.

### 5. QA Prompt for a Paper Discussion

#### 5.1. Discussion Question

How this new method could be beneficial in everyday applications?

#### 5.2. Your Answer

The technique proposed in the paper is a way for computers to learn about videos without any human-provided labels or annotations. Imagine it as a computer watching lots of videos and learning from them, just like we learn from watching. This can be super useful in everyday applications.

For example, think about video recommendation systems like those used by YouTube or Netflix. They recommend videos based on what you've watched before. With "VideoMAE," these systems could get even better at understanding what you like because they can learn more from the videos themselves, without needing humans to categorize or label them.

Another application could be in security cameras. "VideoMAE" could help these cameras recognize unusual or suspicious activities on their own, without needing someone to watch the footage all the time.