# Motion Guided Token Compression for Efficient Masked Video Modeling

**Yukun Feng    Yangming Shi    Fengze Liu    Tan Yan**

ByteDance
{yukun.feng1, shiyangming.best, fengze.liu, tan.yan}@bytedance.com

## Abstract

Recent developments in Transformers have achieved notable strides in enhancing video comprehension. Nonetheless, the $O(N^2)$ computation complexity associated with attention mechanisms presents substantial computational hurdles when dealing with the high dimensionality of videos. This challenge becomes particularly pronounced when striving to increase the frames per second (FPS) to enhance the motion capturing capabilities. Such a pursuit is likely to introduce redundancy and exacerbate the existing computational limitations. In this paper, we initiate by showcasing the enhanced performance achieved through an escalation in the FPS rate. Additionally, we present a novel approach, Motion Guided Token Compression (MGTC), to empower Transformer models to utilize a smaller yet more representative set of tokens for comprehensive video representation. Consequently, this yields substantial reductions in computational burden and remains seamlessly adaptable to increased FPS rates. Specifically, we draw inspiration from video compression algorithms and scrutinize the variance between patches in consecutive video frames across the temporal dimension. The tokens exhibiting a disparity below a predetermined threshold are then masked. Notably, this masking strategy effectively addresses video redundancy while conserving essential information. Our experiments, conducted on widely examined video recognition datasets, Kinetics-400, UCF101 and HMDB51, demonstrate that elevating the FPS rate results in a significant top-1 accuracy score improvement of over 1.6, 1.6 and 4.0. By implementing MGTC with the masking ratio of 25%, we further augment accuracy by 0.1 and simultaneously reduce computational costs by over 31% on Kinetics-400. Even within a fixed computational budget, higher FPS rates paired with MGTC sustain performance gains when compared to lower FPS settings.

## 1    Introduction

In recent years, Transformer-based methods have achieved significant improvement in video understanding(Girdhar et al. 2019; Sharir, Noy, and Zelnik-Manor 2021; Neimark et al. 2021; Liu et al. 2022; Khan et al. 2022; Wei et al. 2022; Ryali et al. 2023; Feichtenhofer et al. 2022; Wang et al. 2022b; Selva et al. 2023). The attention mechanism has been demonstrated the effectiveness in building the de-
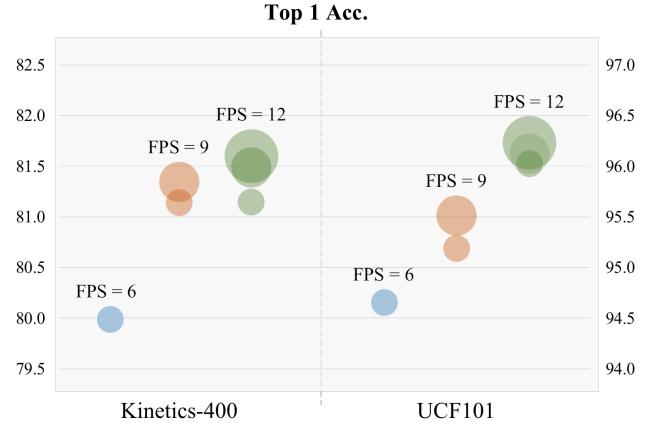
Figure 1: **Overview of the Top-1 Accuracy under different FPS on Kinetics-400 and UCF101.** The relative circle size symbolizes the comparative extent of the computational capacity, which is controlled by the masking ratio and FPS.

pendency between visual tokens(Chen et al. 2020; Dosovitskiy et al. 2021; Liu et al. 2021a). However, it naturally suffers from a $O(N^2)$ computational cost with respect to the input sequence length, which further limits the FPS when we feed video frames. To tackle this issue, researches try to decompose the spatial and temporal dimension and use two attention modules separately (Arnab et al. 2021; Bertasius, Wang, and Torresani 2021). Such methods can largely reduce the memory cost but it becomes challenging to learn the joint spatial-temporal relationships between video frames, and downgrades the performance eventually.

Due to such a limitation, most works (Selva et al. 2023) fail to incorporate a higher FPS and they normally evaluate the model under a fixed rate. However, as introduced in (Mackin, Zhang, and Bull 2019), a higher FPS allows for more temporal information to be captured, not only by the humans, but also benefit the neural video models with a better representation of video actions. This is particularly crucial for tasks that need subtle movements and temporal patterns, such as the action recognition (Kay et al. 2017; Soomro, Zamir, and Shah 2012; Goyal et al. 2017; Ulhaq et al. 2022). Moreover, a higher FPS is able to reduce

the temporal aliasing effect, known as the "strobe effect". When the fast-moving objects in low-FPS is likely to appear choppy or blurred, which confuses the video models to understand the content (Borsato, Aluani, and Morimoto 2015). Increasing the FPS allows the model to reduce the impact of temporal aliasing, which we believe would result in the better performance for neural video models as well. Yet, it is acknowledged that increased FPS rates are likely to introduce additional redundancy, and increase the computational cost, particularly given the current attention mechanisms.

In response to the inherent redundancy in videos, video compression algorithms (Chen, Kao, and Lin 2006) have become a widespread subject of study and application in contemporary video encoding. These algorithms primarily concentrate on compressing motion data while maintaining the fidelity of the reconstructed video. The key frames and motion vectors extracted from a compressed video format have demonstrated their utility across various visual tasks, as illustrated in (Zheng et al. 2023). Nonetheless, the potential sparsity of these motion vectors is not fully harnessed, leading to persistent concerns regarding redundancy and computational overhead when integrating compressed videos into computer vision frameworks.

Recently, the application of masked video modeling has demonstrated its effectiveness in the domain of representation learning (Tong et al. 2022; Feichtenhofer et al. 2022; Wang et al. 2023a). This approach involves a self-supervised task wherein the entire video is reconstructed using a restricted subset of randomly chosen tokens, and yields reasonable reconstruction quality, indicating that spatial-temporal redundancy can be addressed using only a fraction of the video tokens. However, it's important to note that this masking is typically implemented solely during the pre-training phase and doesn't alleviate the computational constraints during the fine-tuning stage. Additionally, the strategy of random token selection still imposes limitations on performance, as each video token shares the same probability of being masked, regardless of its level of informativeness. Consequently, essential tokens representing crucial video information might end up being masked.

In an effort to enhance the masking approach, we aim to maintain a high FPS rate while diminishing redundancy and adhering to a predefined computational budget. Our inspiration predominantly stems from video compression algorithms, propelling us to introduce Motion Guided Token Compression (MGTC) to address these challenges. Contrary to the random token masking technique, our approach involves uniformly segmenting the entire video into non-overlapping patches (or cubes, when the temporal block range exceeds 1). For each token, we compute the disparity between itself and the token at the corresponding spatial position in subsequent frames. This computation signifies the variance in motion across the temporal dimension for the current token. Subsequently, we retain tokens with a motion difference surpassing a predefined threshold. With this method, MGTC guarantees the retention of more informative tokens while effectively discarding redundant tokens through this lightweight token-difference mechanism. Importantly, MGTC remains compatible with higher FPS rates,

as its masking strategy ensures a consistent number of input tokens. This approach can be seamlessly integrated into either the evaluation or training phase, effectively alleviating computational constraints. As demonstrated in Figure 1, we present compelling evidence of the consistent performance enhancement achieved through higher FPS rates within a fixed computational budget, across three extensively studied video datasets, Kinetics-400, UCF101 and HMDB51. Additionally, we observe that implementing MGTC with a designated masking ratio further amplifies accuracy while reducing computational costs. To summarize our contributions:

- We demonstrate the performance gain when increasing the FPS rate. Given a fixed computational budget, a higher FPS is even better with masking 50% tokens.

- We propose the Motion Guided Token Compression strategy, consisting of lightweight motion-guided masking, which is capable of removing spatial-temporal redundancy as well as keeping the action movement, which is able to address the the computational limitation during inference. Also, MGTC is able to enhance video representation by applying the masking during training, and further pushes the performance gain.

- We find the MGTC is better than other masking methods under different masking ratio. It is even better than using all video tokens when we mask 10-20% tokens, demonstrating its necessary in removing the video redundancy.

## 2  Related Work

**Video Action Recognition**  Video action recognition is a representative task in video understanding, and it has seen significant advancements in recent years. Early research in this field mainly focused on CNN-based approaches, which address the video temporal information with a two-stream CNN network (Fan et al. 2016; Feichtenhofer et al. 2019a) or the 3D-CNN directly (Cheron, Laptev, and Schmid 2015; Ye et al. 2015; Diba et al. 2017; Hou, Chen, and Shah 2017; Babaee, Dinh, and Rigoll 2018; Ullah et al. 2018; Hegde et al. 2018; Tu et al. 2018; Lin, Gan, and Han 2019; Yao, Lei, and Zhong 2019; Sudhakaran, Escalera, and Lanz 2020; Li et al. 2020; Duta et al. 2021).

However, inspired by the success of Transformer in natural language processing (Vaswani et al. 2017; Devlin et al. 2019) and image modeling (Dosovitskiy et al. 2021; Liu et al. 2021b) tasks, recent studies (Arnab et al. 2021; Liu et al. 2022; Selva et al. 2023; Bertasius, Wang, and Torresani 2021; Patrick et al. 2021) have explored the application of the Transformer architecture to video action recognition. For instance, ViViT (Arnab et al. 2021) and Timesformer (Bertasius, Wang, and Torresani 2021) have explored factorizing the spatial and temporal dimensions of the input video. VideoSwin (Liu et al. 2022) and MViT (Fan et al. 2021; Li et al. 2022c; Ryali et al. 2023) have proposed hierarchical structures to enable video representation learning and introduce inductive bias to vision transformers. Uniformer and its improvements (Li et al. 2022a,b) have designed a hybrid backbone by integrating 3D-CNN to transformers, combining the advantages of both. These recent advancements along with other works (Yang et al. 2022; Chen et al. 2022;

Ulhaq et al. 2022; Selva et al. 2023) in video action recognition have made great progress and shown promising results.

**Masked Vision Modeling** The central objective of masked vision modeling is to develop proficient representations of visual data by reconstructing data from deliberately corrupted sources. This concept draws inspiration from the Masked Language Modeling task pioneered by BERT(Devlin et al. 2019; Liu et al. 2019). To achieve the best outcomes, it is essential to design a well-constructed reconstruction task. The reconstructed targets can encompass either the original pixel-level data or derived features, both of which have been extensively explored in recent research (Chen et al. 2020; He et al. 2022; Bao et al. 2022; Peng et al. 2022; Xie et al. 2022; Wang et al. 2022a; Girdhar et al. 2023; Wang et al. 2023b; Bandara et al. 2023; Sun et al. 2023). The choice of the masking strategy in masked video modeling is of utmost importance. One effective strategy, called "Tube masking," was introduced by VideoMAE(Tong et al. 2022), and it has proven superior to other strategies like frame masking and random masking. Tube masking prevents information leakage during pre-training, contributing to its success. Additionally, MAR (Motion-AR) (Qing et al. 2022) introduced the concept of Cell Running masking to enhance fine-tuning efficiency, and this concept was further utilized in VideoMAE-v2 (Wang et al. 2023a) for large-scale video pre-training. Combining these established techniques with advanced video compression algorithms, we introduce a groundbreaking motion-guided masking approach for tokenizing compressed videos. This innovative approach promises to improve the efficiency and effectiveness of video tokenization while also minimizing redundancy in the videos and reducing computational demands during the fine-tuning and inference process.

**Video Compression** Video compression is a widely researched area, and numerous works have contributed to the development of efficient and effective video encoding methods (Girod et al. 2005; Adami, Signoroni, and Leonardi 2007; Kumar 2019). One of the seminal works in video compression is the H.264/AVC standard (Sullivan and Wiegand 2005; Zhao and Liang 2006), which introduced advanced video coding techniques such as motion compensation, transform coding, and entropy coding. This standard significantly improved video compression efficiency compared to its predecessors. Another important contribution is the High Efficiency Video Coding (HEVC) standard (Sullivan et al. 2012), also known as H.265, which further enhanced compression efficiency by incorporating advanced coding tools like larger block sizes, improved motion compensation, and more sophisticated entropy coding.

In recent years, there has been a growing interest in deep learning-based video compression methods (Ma et al. 2019; Li, Li, and Lu 2021; Yang et al. 2021). These approaches leverage neural networks to learn and exploit temporal and spatial redundancies in videos, leading to improved compression performance. Additionally, research efforts have focused on exploring emerging video compression techniques such as perceptual video coding, which takes into account human visual perception to allocate bits more ef-

ficiently. Overall, the related works in video compression have made significant strides in achieving higher compression ratios while maintaining acceptable video quality, enabling efficient storage and transmission of video content in various applications.

# 3 Methodology

In order to ensure the compatibility of the video model with an higher FPS rate while concurrently managing feature redundancy within a constrained computational budget, we introduce a straightforward yet efficient masking strategy named Motion Guided Token Compression (MGTC). This strategy effectively handles the redundancy by selectively retaining only the informative video patches found between frames. Notably, MGTC can be seamlessly integrated during the inference stage, maintaining the original computational cost even when the FPS is enhanced. Furthermore, the application of MGTC during the training phase could further accelerate the training and decrease the computational cost.

Figure 2 shows an example of the comparison of MGTC and other two masking strategies, under a masking ratio of 50%. It becomes evident that MGTC not only retains a greater amount of action dynamics but also effectively mitigates feature redundancy across the temporal dimension, especially in a higher FPS setting.

## 3.1 Motion Guided Token Compression

The conventional video compression approach, such as H264 (Chen, Kao, and Lin 2006), compresses video content by encoding pixel data through the identification and comparison of variations between predicted and observed pixels. This process allows for the discrimination between essential and redundant pixel information.

Taking inspiration from the encoding principles of H264, MGTC employs a similar concept to eliminate duplicate video patches between consecutive frames. This is achieved through an analysis of pixel disparities across the temporal dimension, encompassing two distinct steps: Sub-block Division and Block Masking. Figure 3 provides an illustrative instance that shows the MGTC tokenization process, delineating how the frame is partitioned into patches and specifying the patches that are retained.

**Sub-block Division** MGTC first divides the video into multiple non-overlapping sub-blocks. Given a video with dimensions of $T \times H \times W$, we use a cube with dimensions of $c \times p1 \times p2$ to tokenize it into a sequence of $L = \frac{T}{c} \times \frac{H}{p1} \times \frac{W}{p2}$ blocks. By treating the video as a series of sub-blocks, this sequence effectively retains all the details present in the original video. This approach not only improves the ability to capture subtle distinctions in subsequent stages but also seamlessly address the limitation of transformer-based models, thereby creating a practical and feasible input format suitable for diverse applications.

**Block Masking** Since a significant portion of video frames consists of duplicate content, the primary distinction typically revolves around the motion dynamics. MGTC efficiently identifies changes in action by examining pixel
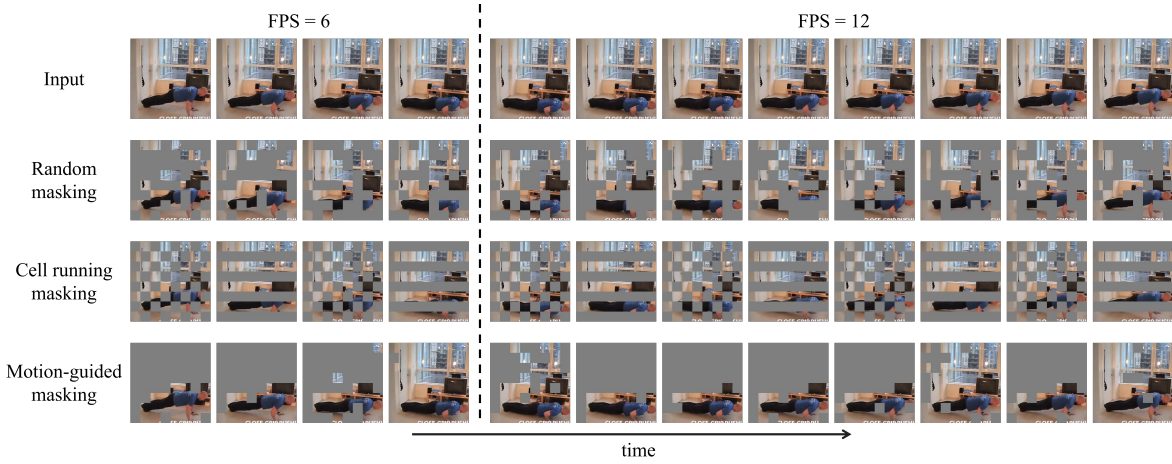
Figure 2: **Comparison between different masking methods, under various FPS rate.** MGTC is able to capture the action movement, and remove the redundant information, especially in higher FPS rate. Here we use a masking ratio of 50%.
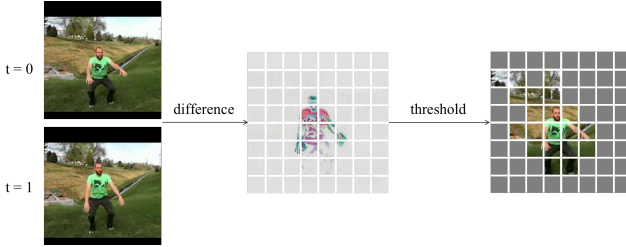


Figure 3: **Workflow of motion-guided masking.** Patch differences are calculated for further masking.

differences between consecutive blocks within the temporal dimension. Specifically, blocks that exhibit minor or no changes in succession are considered redundant and, consequently, are masked from the input sequence.

As we define the above sequence of blocks with a length of $L$, before flattening it into a 1D sequence, we keep its spatial and temporal dimension as $\{\frac{T}{c}\} \times \{\frac{H}{p1} \times \frac{W}{p2}\}$. Then MGTC compares the pixel difference $D$ with the adjacent available blocks along the temporal dimension. If $D$ is below a threshold $\lambda$, the blocks will be masked, described as:

$$Mask_j^i = True \text{ if } D_j^i < \lambda \text{ else } False,$$
$$D_j^i = MSE(C_j^i, \ C_j^{i+1})$$

where the $i$ and $j$ refer to the index of the temporal and spatial dimension respectively. Hence, blocks exhibiting notable pixel alterations are grouped together to effectively represent the motion patterns within the input video. For the purpose of representing video appearances, a temporal index is randomly chosen during the training process, similar to a key frame in H264 (Chen, Kao, and Lin 2006). All video patches originating from this selected index are retained.

Within MGTC, the hyperparameter $\lambda$ plays a pivotal role in determining the quantity of cubes necessitating masking. Given the inherent diversity in videos, utilizing a uniform threshold across the entire dataset is unfeasible. Instead, we dynamically compute the threshold for each specific video based on its unique attributes. Our approach assumes that a certain proportion of each video constitutes redundancy. Consequently, we set $\lambda$ to match the cube difference value at the upper $N$ percent of all cubes present within that video. In the forthcoming analysis section, we will delve into the exploration of varying values for $N$ to assess their impact.

## 3.2 Training and Evaluation

MGTC serves as a versatile plug-and-play enhancement, capable of seamless integration with a diverse range of transformer-based models for video comprehension, whether during training or evaluation. It's essential to acknowledge that earlier researchers have identified the potency of joint time-space attention across all video cubes. However, this approach introduces a computationally intensive complexity of $N^2$. Consequently, we opt to utilize VideoMAE as a baseline model, aiming to validate the efficiency and effectiveness of MGTC in this context.

**Evaluation with MGTC** Rather than incorporate MGTC into the training process, we can opt for a simpler yet effective approach. We utilize the representations that have been thoroughly trained with all video tokens during fine-tuning. In the inference stage, we selectively employ informative tokens that were selected by MGTC. This method proves to be straightforward, which leverages the advantages of training with the full token set while significantly reducing the number of tokens used during inference. It also alleviates the inference bottleneck while still preserving the essential information within the video.

**Training with MGTC** We adhere to the training configuration outlined in VideoMAE (Tong et al. 2022), which consists of an initial pre-training phase involving VideoMAE reconstruction, followed by fine-tuning for classification tasks on downstream objectives. In this study, we do not incorporate MGTC during the pre-training stage. Instead, we use the same setup as the original VideoMAE. The primary dis-

tinction lies in our utilization of a higher FPS during the pre-training phase, leading to an increase in computational demands. During the fine-tuning phase, we explore two approaches: either using all video tokens, as per the original setting, or applying MGTC to retain informative video cubes that can represent the entire video for training purposes. We anticipate that integrating further training with MGTC will amplify the representation quality of the chosen video patches, concurrently expediting the training process.

# 4 Experiments

## 4.1 Datasets

Our experimentation encompasses two extensively examined video action recognition datasets: Kinetics-400 (Kinetics-400) (Kay et al. 2017), UCF101 (Soomro, Zamir, and Shah 2012) and HMDB51 (Kuehne et al. 2011). Kinetics-400 offers a substantial scale with around 240k training videos and 20k validation videos, each spanning 10 seconds in duration. This dataset encompasses a comprehensive range of 400 distinct classes. In contrast, UCF101 and HMDB51 are more compact datasets, featuring 9.5k/3.5k and 3.5k/1.5k training as well as validation videos.
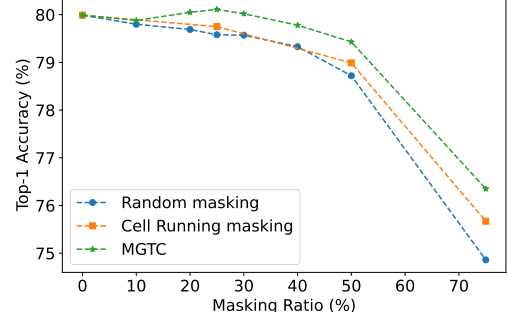
## 4.2 Settings

In line with the pre-training and fine-tuning configurations established in VideoMAE (Tong et al. 2022), we employ the masking strategy outlined in Figure 3 for either the fine-tuning or inference stages. For the remaining training and evaluation hyperparameters, we preserve the same values as those used in the original VideoMAE.
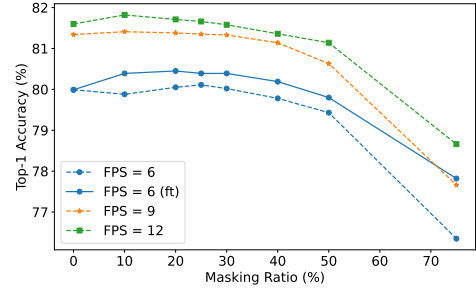
## 4.3 Ablation Study

We have demonstrated the performance improvement resulting from a higher FPS rate, as highlighted in Figure 1. Given that incorporating additional video frames significantly escalates computational demands, we have compared three masking strategies aimed at mitigating this limitation. Remarkably, MGTC consistently outperforms both the Random and Cell Running masking (Qing et al. 2022). Also, surprisingly, the model achieves higher accuracy even when employing a relatively small percentage of masking. Ablation experiments are conducted with the Kinetics-400 dataset, in which we pre-train the model for 800 epochs while maintaining all other parameters the same.

**Masking Methods**  MGTC effectively handles redundancy by masking duplicates while retaining essential tokens, a capability not only found in other methods such as simple random masking and the Cell Running masking (Qing et al. 2022), but also surpassing them in performance, as illustrated in Figure 4a. This figure clearly demonstrates that MGTC consistently outperforms the other two methods across various masking ratios, highlighting its ability to preserve more representative tokens. Furthermore, when applying a masking ratio of 10%-20%, MGTC achieves superior accuracy scores compared to using all patch tokens, in contrast to the declining scores observed with the other two methods as the masking ratio increases. This observation serves as further evidence that MGTC excels at reducing



(a) Comparison of masking methods under different masking ratios. Experiments are made with 6-FPS setting. MGTC is consistently better than other Random and Cell Running masking. MGTC with 10% masking even outperforms using all video tokens.



(b) Comparison of MGTC in different FPS. A higher FPS brings up performance gain under each masking ratio. "(ft)" refers to training with MGTC.

video redundancy, leading to enhanced performance without significant additional computational cost.

**Masking under Different FPS**  In Figure 4b, we present a comparison of MGTC performance across varying FPS settings. As the FPS increases, the performance score demonstrates an upward trend and a noteworthy enhancement is observed when transitioning from 6 to 9 FPS, although the relative benefit diminishes when further increasing the FPS to 12. In settings with higher FPS values, the model reaches its peak relative performance gain with a masking ratio of 10%, whereas this ratio is 25% for lower FPS scenarios. This consistency suggests that a higher FPS captures more action-related information, resulting in a lower relative redundancy percentage. Additionally, it's interesting to note that the relative performance drop in higher FPS settings diminishes when utilizing a masking ratio of 50%. This indicates that a higher FPS introduces more redundancy, which, however, is effectively addressed by the MGTC, exerting minimal influence on overall performance.

**Masking during Training**  By incorporating MGTC during the training phase, the outcomes depicted in Figure 4b underscore the effectiveness of training the selected video tokens. Training with MGTC enables the model to improve the representation of preserved tokens, resulting in an enhanced video representation and a more substantial performance boost.

| Model | Backbone | Pre-train | FPS | GFLOPs (G) | Param (M) | Top-1 (%) | Top-5 (%) |
|---|---|---|---|---|---|---|---|
| SlowFast | Res101+NL | IN21K | 12+3 | $359 \times 10 \times 3$ | 60 | 79.8 | 93.9 |
| Timesformer | ViT-L | IN21K | 6 | $8353 \times 1 \times 3$ | 430 | 80.7 | 94.7 |
| ViViT FE | ViT-L | IN21k | 24 | $3980 \times 1 \times 3$ | 430 | 81.7 | 93.8 |
| Motionformer | ViT-L | IN21K | 12 | $1185 \times 10 \times 3$ | 382 | 80.2 | 94.8 |
| VideoSwin | Swin-L | IN21K | 12 | $604 \times 4 \times 3$ | 197 | 83.1 | 95.9 |
| MViTv1 | MViTv1-B | IN21K | 12 | $170 \times 5 \times 1$ | 37 | 80.2 | 94.3 |
| BEVT | Swin-B | IN-1K+DALLE | 12 | $282 \times 4 \times 3$ | 88 | 80.6 | N/A |
| OmniMAE | ViT-B | IN-1K+Kinetics-400 | 12 | $180 \times 5 \times 3$ | 87 | 80.6 | N/A |
| MaskFeat | MViTv1-L | Kinetics-400 | 6 | $377 \times 10 \times 1$ | 218 | 84.3 | 96.3 |
| MME | ViT-B | Kinetics-400 | 6 | $180 \times 7 \times 3$ | 87 | 81.8 | N/A |
| $Ada$MAE | ViT-B | Kinetics-400 | 6 | $180 \times 7 \times 3$ | 87 | 81.7 | 95.2 |
| $MAR_{\rho=50\%}$ | ViT-B | Kinetics-400 | 6 | $86 \times 5 \times 3$ | 94 | 81.0 | 94.4 |
| $MAR_{\rho=50\%}$ | ViT-L | Kinetics-400 | 6 | $276 \times 5 \times 3$ | 311 | 85.3 | 96.3 |
| $VideoMAE_{e=800}$ | ViT-B | Kinetics-400 | 6 | $180 \times 5 \times 3$ | 87 | 80.0 | 94.4 |
| $VideoMAE_{e=800}$ | ViT-B | Kinetics-400 | 9 | $240 \times 5 \times 3$ | 87 | 81.3 | 94.9 |
| $VideoMAE_{e=800}$ | ViT-B | Kinetics-400 | 12 | $451 \times 5 \times 3$ | 87 | 81.6 | 94.9 |
| $\mathbf{MGTC}_{\rho=25\%,e=800}$ | ViT-B | Kinetics-400 | 6 | $127 \times 5 \times 3$ | 87 | 80.4 | 94.5 |
| $\mathbf{MGTC}_{\rho=10\%,e=800}$ | ViT-B | Kinetics-400 | 9 | $210 \times 5 \times 3$ | 87 | 81.4 | 94.8 |
| $\mathbf{MGTC}_{\rho=10\%,e=800}$ | ViT-B | Kinetics-400 | 12 | $392 \times 5 \times 3$ | 87 | 81.8 | 94.9 |
| $VideoMAE_{e=1600}$ | ViT-B | Kinetics-400 | 6 | $180 \times 5 \times 3$ | 87 | 81.5 | 95.0 |
| $VideoMAE_{e=1600}$ | ViT-B | Kinetics-400 | 12 | $451 \times 5 \times 3$ | 87 | 81.8 | 95.0 |
| $\mathbf{MGTC}_{\rho=25\%,e=1600}$ | ViT-B | Kinetics-400 | 6 | $127 \times 5 \times 3$ | 87 | 81.6 | 95.0 |
| $\mathbf{MGTC}_{\rho=10\%,e=1600}$ | ViT-B | Kinetics-400 | 12 | $392 \times 5 \times 3$ | 87 | 82.0 | 95.1 |
| $VideoMAE_{e=1600}$ | ViT-L | Kinetics-400 | 6 | $597 \times 5 \times 3$ | 305 | 85.2 | 96.8 |
| $VideoMAE_{e=1600}$ | ViT-L | Kinetics-400 | 12 | $1436 \times 5 \times 3$ | 305 | 85.4 | 97.0 |
| $\mathbf{MGTC}_{\rho=50\%,e=1600}$ | ViT-L | Kinetics-400 | 6 | $269 \times 5 \times 3$ | 305 | 85.3 | 97.0 |
| $\mathbf{MGTC}_{\rho=25\%,e=1600}$ | ViT-L | Kinetics-400 | 12 | $988 \times 5 \times 3$ | 305 | 85.5 | 96.8 |

(a) Comparison with state-of-the-arts on Kinetics-400.

| Model | Backbone | Extra data | FPS | GFLOPs (G) | Param (M) | UCF101 Top-1(%) | HMDB51 Top-1(%) |
|---|---|---|---|---|---|---|---|
| XDC | R(2+1)D | IG65M | 12 | N/A | 15 | 94.2 | 67.1 |
| GDT | R(2+1)D | IG65M | 12 | N/A | 15 | 95.2 | 72.8 |
| CVRL | Res50 | Kinetics-400 | 12 | N/A | 32 | 92.9 | 67.9 |
| $CORP_f$ | Res50 | Kinetics-400 | 3 | N/A | 32 | 87.3 | 68.0 |
| $\rho$BYOL | Res50 | Kinetics-400 | 3 | N/A | 32 | 94.2 | 72.1 |
| MME | ViT-B | Kinetics-400 | $6^*$ | $180 \times 5 \times 3$ | 87 | 96.5 | 78.0 |
| $VideoMAE_{e=800}$ | ViT-B | Kinetics-400 | $6^*$ | $180 \times 5 \times 3$ | 87 | 94.6 | 70.1 |
| $VideoMAE_{e=800}$ | ViT-B | Kinetics-400 | 9 | $240 \times 5 \times 3$ | 87 | 95.5 | N/A |
| $VideoMAE_{e=800}$ | ViT-B | Kinetics-400 | $12^*$ | $451 \times 5 \times 3$ | 87 | 96.2 | 74.1 |
| $\mathbf{MGTC}_{\rho=10\%,e=800}$ | ViT-B | Kinetics-400 | $6^*$ | $159 \times 5 \times 3$ | 87 | 94.7 | 70.7 |
| $\mathbf{MGTC}_{\rho=25\%,e=800}$ | ViT-B | Kinetics-400 | 9 | $167 \times 5 \times 3$ | 87 | 95.2 | N/A |
| $\mathbf{MGTC}_{\rho=50\%,e=800}$ | ViT-B | Kinetics-400 | $12^*$ | $182 \times 5 \times 3$ | 87 | 96.0 | 73.3 |
| $VideoMAE_{e=1600}$ | ViT-B | Kinetics-400 | $6^*$ | $180 \times 5 \times 3$ | 87 | 95.4 | 72.2 |
| $VideoMAE_{e=1600}$ | ViT-B | Kinetics-400 | $12^*$ | $451 \times 5 \times 3$ | 87 | 96.2 | 74.1 |
| $\mathbf{MGTC}_{\rho=50\%,e=1600}$ | ViT-B | Kinetics-400 | $6^*$ | $80 \times 5 \times 3$ | 87 | 95.5 | 71.9 |
| $\mathbf{MGTC}_{\rho=25\%,e=1600}$ | ViT-B | Kinetics-400 | $12^*$ | $306 \times 5 \times 3$ | 87 | 96.3 | 74.2 |

(b) Comparison with the state-of-the-arts on UCF101 and HMDB51. The superscript $^*$ indicates that it is multiplied by 2 on the HMDB dataset.

Table 1: **System-level comparisons on Kinetics-400, UCF101 and HMDB51 dataset.** The GFLOPs refers to 'FLOPs $\times$ Clips $\times$ Crops' The subscript $\rho$ is the mask ratio during inference, while $e$ is the pre-training epoch.

## 4.4 Main Results

We compare different FPS settings with either the MGTC or the no-masking strategy in Kinetics-400, UCF101 and HMDB51, the results are shown in Table 1. The SOTA methods we use for comparison are SlowFast (Feichtenhofer et al. 2019b), Timesformer (Bulat et al. 2021), ViViT FE (Arnab et al. 2021), MotionFormer (Patrick et al. 2021), VideoSwin (Liu et al. 2021c), MViTv1 (Fan et al. 2021), BEVT (Wang et al. 2022a), OmniMAE (Girdhar et al. 2023), MaskFeat (Wei et al. 2022), MME (Sun et al. 2023), $ada$MAE (Bandara et al. 2023), MAR (Qing et al. 2022), VideoMAE (Tong et al. 2022), XDC (Alwassel et al. 2020), GDT (Patrick et al. 2020), CVRL (Qian et al. 2021), CORP$_f$ (Hu et al. 2021), $\rho$BYOL (Feichtenhofer et al. 2021).

It is noticeable that there is a direct relationship between the FPS and the accuracy scores. When the FPS is increased, for all datasets, we consistently observe an improvement in the top-1 accuracy score, from 80.0%, 94.6%, 70.1% to 81.6%, 96.2%, 74.1%, which suggests that optimizing FPS can bring up more motion information, leading to a positive impact on the performance.

Also, when we use a mask ratio of 25% for 6 FPS and 10% for the higher FPS setting, there is an improvement up to 0.2 accuracy score. Slight masking with the MGTC not only further enhances the performance, but also does so in a more resource-efficient manner, which could be interpreted as removing the video redundancy decrease the noise of input tokens and further improve the accuracy eventually.

Another noteworthy finding from the table is that even when we operate within a fixed budget, the higher FPS settings with MGTC consistently outperform the lower FPS settings. This result demonstrates that investing resources in optimizing FPS with our approach is a worthwhile strategy, as it yields better results within the same budgetary limits, making the most of the available resources.

To ensure the scaling ability of the MGTC, we conducted extensive experiments using pre-training on large models, including ViT-Large, over an extended training period of 1600 epochs. MGTC consistently produces positive results with a 0.2 gain from higher FPS and an additional 0.1 gain from MGTC. In this demanding context, it demonstrates its reliability and potential for broader applicability in other various scenarios.

## 4.5 Discussions

**Pixel-Residual Distribution** In this paper, we introduce a video redundancy elimination strategy built upon patch similarity. For the sake of simplicity, we employ the frame difference as a measure of similarity between patches. Our experiments have demonstrated that discarding redundant information not only can diminish computational demands, but also enhance model performance particularly when the mask ratio is lower.

To further underscore the prevalence of excessive redundant data within videos, we have highlighted the frame difference distribution in Figure 5. It can be observed that approximately 25% of frame differences within the Kinetics-400 dataset are near zero, and astonishingly, this proportion
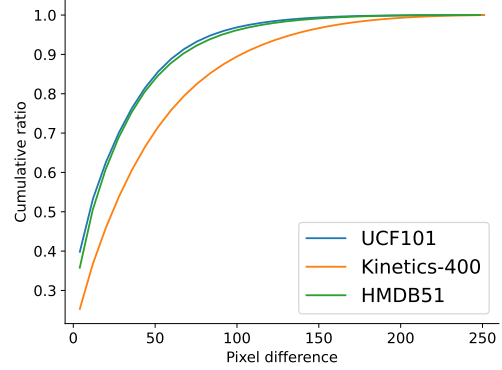


Figure 5: Pixel-Residual Distributions under 12 FPS.

escalates to 40% in the case of UCF101. Nevertheless, such redundant data can place an undue burden on the model, hence, their removal can contribute to improved modelling performance. For instance, based on the same model applied to the Kinetics-400 dataset, increasing the mask ratio from 0% to 40% improved the accuracy of the model from 80.0% to a slightly higher 80.4%.

When adopting a higher FPS and more frames, while also using MGTC to mask certain patches (for instance, 12-FPS, mask ratio=50%), as opposed to when using 6-FPS without masking, we essentially eliminate low information density data (such as background details) in favor of retaining more high information density data (such as motion content). The result is a substantial boost in overall performance.

**Computational Complexity** We have showcased the efficacy of capturing additional motions through a higher FPS rate. Nevertheless, it is essential to acknowledge that this approach also introduces redundancy and imposes supplementary computational burdens. In this context, we emphasize that MGTC serves to alleviate the augmented computational load by applying masking to a specific percentage of video tokens, as illustrated in Table 1. Consequently, the computational expense of employing higher FPS rates with MGTC aligns with that of lower FPS rates, while delivering substantial performance improvements. In essence, MGTC facilitates the maintenance of superior performance within a fixed computational budget for higher FPS settings.

## 5 Conclusion

This paper demonstrates the advantages gained from increasing the FPS rate. Additionally, we introduce MGTC as a means to amplify video token representation by eliminating redundant patches, simultaneously reducing computational expenses within a predefined budget. Our primary aim is to emphasize the significance of elevating FPS rates in video action recognition, with MGTC representing just one among several potential solutions to tackle computational limitations. We encourage future research to dive into the merits of higher FPS settings and to explore more intuitive approaches for alleviating computational constraints.

# References

Adami, N.; Signoroni, A.; and Leonardi, R. 2007. State-of-the-Art and Trends in Scalable Video Compression With Wavelet-Based Approaches. *IEEE Transactions on Circuits and Systems for Video Technology*, 17(9): 1238–1255.

Alwassel, H.; Mahajan, D.; Korbar, B.; Torresani, L.; Ghanem, B.; and Tran, D. 2020. Self-supervised learning by cross-modal audio-video clustering. *Advances in Neural Information Processing Systems*, 33: 9758–9770.

Arnab, A.; Dehghani, M.; Heigold, G.; Sun, C.; Lučić, M.; and Schmid, C. 2021. ViViT: A Video Vision Transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 6836–6846.

Babaee, M.; Dinh, D. T.; and Rigoll, G. 2018. A deep convolutional neural network for video sequence background subtraction. *Pattern Recognition*, 76: 635–649.

Bandara, W. G. C.; Patel, N.; Gholami, A.; Nikkhah, M.; Agrawal, M.; and Patel, V. M. 2023. AdaMAE: Adaptive Masking for Efficient Spatiotemporal Learning with Masked Autoencoders. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 14507–14517.

Bao, H.; Dong, L.; Piao, S.; and Wei, F. 2022. BEiT: BERT Pre-Training of Image Transformers. arXiv:2106.08254.

Bertasius, G.; Wang, H.; and Torresani, L. 2021. Is Space-Time Attention All You Need for Video Understanding? arXiv:2102.05095.

Borsato, F. H.; Aluani, F. O.; and Morimoto, C. H. 2015. A Fast and Accurate Eye Tracker Using Stroboscopic Differential Lighting. In *2015 IEEE International Conference on Computer Vision Workshop (ICCVW)*, 502–510.

Bulat, A.; Perez-Rua, J.-M.; Sudhakaran, S.; Martinez, B.; and Tzimiropoulos, G. 2021. Space-time Mixing Attention for Video Transformer. arXiv:2106.05968.

Chen, J.-W.; Kao, C.-Y.; and Lin, Y.-L. 2006. Introduction to H. 264 advanced video coding. In *Proceedings of the 2006 Asia and South Pacific Design Automation Conference*, 736–741.

Chen, M.; Radford, A.; Child, R.; Wu, J.; Jun, H.; Luan, D.; and Sutskever, I. 2020. Generative pretraining from pixels. In *International conference on machine learning*, 1691–1703. PMLR.

Chen, S.; Ge, C.; Tong, Z.; Wang, J.; Song, Y.; Wang, J.; and Luo, P. 2022. Adaptformer: Adapting vision transformers for scalable visual recognition. *Advances in Neural Information Processing Systems*, 35: 16664–16678.

Cheron, G.; Laptev, I.; and Schmid, C. 2015. P-CNN: Pose-Based CNN Features for Action Recognition. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*.

Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 4171–4186. Minneapolis, Minnesota: Association for Computational Linguistics.

Diba, A.; Fayyaz, M.; Sharma, V.; Karami, A. H.; Arzani, M. M.; Yousefzadeh, R.; and Gool, L. V. 2017. Temporal 3D ConvNets: New Architecture and Transfer Learning for Video Classification. arXiv:1711.08200.

Dosovitskiy, A.; Beyer, L.; Kolesnikov, A.; Weissenborn, D.; Zhai, X.; Unterthiner, T.; Dehghani, M.; Minderer, M.; Heigold, G.; Gelly, S.; Uszkoreit, J.; and Houlsby, N. 2021. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In *International Conference on Learning Representations*.

Duta, I. C.; Liu, L.; Zhu, F.; and Shao, L. 2021. Improved Residual Networks for Image and Video Recognition. In *2020 25th International Conference on Pattern Recognition (ICPR)*, 9415–9422.

Fan, H.; Xiong, B.; Mangalam, K.; Li, Y.; Yan, Z.; Malik, J.; and Feichtenhofer, C. 2021. Multiscale vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, 6824–6835.

Fan, Y.; Lu, X.; Li, D.; and Liu, Y. 2016. Video-Based Emotion Recognition Using CNN-RNN and C3D Hybrid Networks. In *Proceedings of the 18th ACM International Conference on Multimodal Interaction*, ICMI '16, 445–450. New York, NY, USA: Association for Computing Machinery. ISBN 9781450345569.

Feichtenhofer, C.; Fan, H.; Malik, J.; and He, K. 2019a. SlowFast Networks for Video Recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.

Feichtenhofer, C.; Fan, H.; Malik, J.; and He, K. 2019b. SlowFast Networks for Video Recognition. arXiv:1812.03982.

Feichtenhofer, C.; Fan, H.; Xiong, B.; Girshick, R.; and He, K. 2021. A large-scale study on unsupervised spatiotemporal representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 3299–3309.

Feichtenhofer, C.; Li, Y.; He, K.; et al. 2022. Masked autoencoders as spatiotemporal learners. *Advances in neural information processing systems*, 35: 35946–35958.

Girdhar, R.; Carreira, J.; Doersch, C.; and Zisserman, A. 2019. Video Action Transformer Network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.

Girdhar, R.; El-Nouby, A.; Singh, M.; Alwala, K. V.; Joulin, A.; and Misra, I. 2023. Omnimae: Single model masked pretraining on images and videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 10406–10417.

Girod, B.; Aaron, A.; Rane, S.; and Rebollo-Monedero, D. 2005. Distributed Video Coding. *Proceedings of the IEEE*, 93(1): 71–83.

Goyal, R.; Kahou, S. E.; Michalski, V.; Materzyńska, J.; Westphal, S.; Kim, H.; Haenel, V.; Fruend, I.; Yianilos, P.; Mueller-Freitag, M.; Hoppe, F.; Thurau, C.; Bax, I.; and Memisevic, R. 2017. The "something something" video database for learning and evaluating visual common sense. arXiv:1706.04261.

He, K.; Chen, X.; Xie, S.; Li, Y.; Dollár, P.; and Girshick, R. 2022. Masked Autoencoders Are Scalable Vision Learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 16000–16009.

Hegde, K.; Agrawal, R.; Yao, Y.; and Fletcher, C. W. 2018. Morph: Flexible Acceleration for 3D CNN-Based Video Understanding. In *2018 51st Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, 933–946.

Hou, R.; Chen, C.; and Shah, M. 2017. Tube Convolutional Neural Network (T-CNN) for Action Detection in Videos. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*.

Hu, K.; Shao, J.; Liu, Y.; Raj, B.; Savvides, M.; and Shen, Z. 2021. Contrast and order representations for video self-supervised learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 7939–7949.

Kay, W.; Carreira, J.; Simonyan, K.; Zhang, B.; Hillier, C.; Vijayanarasimhan, S.; Viola, F.; Green, T.; Back, T.; Natsev, P.; Suleyman, M.; and Zisserman, A. 2017. The Kinetics Human Action Video Dataset. arXiv:1705.06950.

Khan, S.; Naseer, M.; Hayat, M.; Zamir, S. W.; Khan, F. S.; and Shah, M. 2022. Transformers in vision: A survey. *ACM computing surveys (CSUR)*, 54(10s): 1–41.

Kuehne, H.; Jhuang, H.; Garrote, E.; Poggio, T.; and Serre, T. 2011. HMDB: a large video database for human motion recognition. In *2011 International conference on computer vision*, 2556–2563. IEEE.

Kumar, D. T. S. 2019. A novel method for HDR video encoding, compression and quality evaluation. *Journal of Innovative Image Processing*, 1(2): 71–80.

Li, D.; Wang, R.; Xie, C.; Liu, L.; Zhang, J.; Li, R.; Wang, F.; Zhou, M.; and Liu, W. 2020. A Recognition Method for Rice Plant Diseases and Pests Video Detection Based on Deep Convolutional Neural Network. *Sensors*, 20(3).

Li, J.; Li, B.; and Lu, Y. 2021. Deep Contextual Video Compression. In Ranzato, M.; Beygelzimer, A.; Dauphin, Y.; Liang, P.; and Vaughan, J. W., eds., *Advances in Neural Information Processing Systems*, volume 34, 18114–18125. Curran Associates, Inc.

Li, K.; Wang, Y.; Gao, P.; Song, G.; Liu, Y.; Li, H.; and Qiao, Y. 2022a. Uniformer: Unified transformer for efficient spatiotemporal representation learning. *arXiv preprint arXiv:2201.04676*.

Li, K.; Wang, Y.; He, Y.; Li, Y.; Wang, Y.; Wang, L.; and Qiao, Y. 2022b. Uniformerv2: Spatiotemporal learning by arming image vits with video uniformer. *arXiv preprint arXiv:2211.09552*.

Li, Y.; Wu, C.-Y.; Fan, H.; Mangalam, K.; Xiong, B.; Malik, J.; and Feichtenhofer, C. 2022c. Mvitv2: Improved multiscale vision transformers for classification and detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 4804–4814.

Lin, J.; Gan, C.; and Han, S. 2019. TSM: Temporal Shift Module for Efficient Video Understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.

Liu, Y.; Ott, M.; Goyal, N.; Du, J.; Joshi, M.; Chen, D.; Levy, O.; Lewis, M.; Zettlemoyer, L.; and Stoyanov, V. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. arXiv:1907.11692.

Liu, Z.; Lin, Y.; Cao, Y.; Hu, H.; Wei, Y.; Zhang, Z.; Lin, S.; and Guo, B. 2021a. Swin Transformer: Hierarchical Vision Transformer using Shifted Windows. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, 9992–10002.

Liu, Z.; Lin, Y.; Cao, Y.; Hu, H.; Wei, Y.; Zhang, Z.; Lin, S.; and Guo, B. 2021b. Swin Transformer: Hierarchical Vision Transformer using Shifted Windows. *CoRR*, abs/2103.14030.

Liu, Z.; Ning, J.; Cao, Y.; Wei, Y.; Zhang, Z.; Lin, S.; and Hu, H. 2021c. Video Swin Transformer. arXiv:2106.13230.

Liu, Z.; Ning, J.; Cao, Y.; Wei, Y.; Zhang, Z.; Lin, S.; and Hu, H. 2022. Video Swin Transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 3202–3211.

Ma, S.; Zhang, X.; Jia, C.; Zhao, Z.; Wang, S.; and Wang, S. 2019. Image and video compression with neural networks: A review. *IEEE Transactions on Circuits and Systems for Video Technology*, 30(6): 1683–1698.

Mackin, A.; Zhang, F.; and Bull, D. R. 2019. A Study of High Frame Rate Video Formats. *IEEE Transactions on Multimedia*, 21(6): 1499–1512.

Neimark, D.; Bar, O.; Zohar, M.; and Asselmann, D. 2021. Video Transformer Network. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops*, 3163–3172.

Patrick, M.; Asano, Y.; Kuznetsova, P.; Fong, R.; Henriques, J. F.; Zweig, G.; and Vedaldi, A. 2020. Multi-modal self-supervision from generalized data transformations. In *International Conference on Learning Representations*.

Patrick, M.; Campbell, D.; Asano, Y.; Misra, I.; Metze, F.; Feichtenhofer, C.; Vedaldi, A.; and Henriques, J. F. 2021. Keeping your eye on the ball: Trajectory attention in video transformers. *Advances in neural information processing systems*, 34: 12493–12506.

Peng, Z.; Dong, L.; Bao, H.; Ye, Q.; and Wei, F. 2022. A Unified View of Masked Image Modeling. arXiv:2210.10615.

Qian, R.; Meng, T.; Gong, B.; Yang, M.-H.; Wang, H.; Belongie, S.; and Cui, Y. 2021. Spatiotemporal contrastive video representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 6964–6974.

Qing, Z.; Zhang, S.; Huang, Z.; Wang, X.; Wang, Y.; Lv, Y.; Gao, C.; and Sang, N. 2022. MAR: Masked Autoencoders for Efficient Action Recognition. arXiv:2207.11660.

Ryali, C.; Hu, Y.-T.; Bolya, D.; Wei, C.; Fan, H.; Huang, P.-Y.; Aggarwal, V.; Chowdhury, A.; Poursaeed, O.; Hoffman, J.; Malik, J.; Li, Y.; and Feichtenhofer, C. 2023. Hiera: A Hierarchical Vision Transformer without the Bells-and-Whistles. *ICML*.

Selva, J.; Johansen, A. S.; Escalera, S.; Nasrollahi, K.; Moeslund, T. B.; and Clapés, A. 2023. Video transformers: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.

Sharir, G.; Noy, A.; and Zelnik-Manor, L. 2021. An Image is Worth 16x16 Words, What is a Video Worth? *ArXiv*, abs/2103.13915.

Soomro, K.; Zamir, A. R.; and Shah, M. 2012. UCF101: A Dataset of 101 Human Actions Classes From Videos in The Wild. arXiv:1212.0402.

Sudhakaran, S.; Escalera, S.; and Lanz, O. 2020. Gate-Shift Networks for Video Action Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.

Sullivan, G.; and Wiegand, T. 2005. Video Compression - From Concepts to the H.264/AVC Standard. *Proceedings of the IEEE*, 93(1): 18–31.

Sullivan, G. J.; Ohm, J.-R.; Han, W.-J.; and Wiegand, T. 2012. Overview of the High Efficiency Video Coding (HEVC) Standard. *IEEE Transactions on Circuits and Systems for Video Technology*, 22(12): 1649–1668.

Sun, X.; Chen, P.; Chen, L.; Li, C.; Li, T. H.; Tan, M.; and Gan, C. 2023. Masked Motion Encoding for Self-Supervised Video Representation Learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2235–2245.

Tong, Z.; Song, Y.; Wang, J.; and Wang, L. 2022. Video-MAE: Masked Autoencoders are Data-Efficient Learners for Self-Supervised Video Pre-Training. In *Advances in Neural Information Processing Systems*.

Tu, Z.; Xie, W.; Qin, Q.; Poppe, R.; Veltkamp, R. C.; Li, B.; and Yuan, J. 2018. Multi-stream CNN: Learning representations based on human-related regions for action recognition. *Pattern Recognition*, 79: 32–43.

Ulhaq, A.; Akhtar, N.; Pogrebna, G.; and Mian, A. 2022. Vision transformers for action recognition: A survey. *arXiv preprint arXiv:2209.05700*.

Ullah, A.; Ahmad, J.; Muhammad, K.; Sajjad, M.; and Baik, S. W. 2018. Action Recognition in Video Sequences using Deep Bi-Directional LSTM With CNN Features. *IEEE Access*, 6: 1155–1166.

Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, L.; and Polosukhin, I. 2017. Attention Is All You Need. *CoRR*, abs/1706.03762.

Wang, L.; Huang, B.; Zhao, Z.; Tong, Z.; He, Y.; Wang, Y.; Wang, Y.; and Qiao, Y. 2023a. VideoMAE V2: Scaling Video Masked Autoencoders with Dual Masking. arXiv:2303.16727.

Wang, R.; Chen, D.; Wu, Z.; Chen, Y.; Dai, X.; Liu, M.; Jiang, Y.-G.; Zhou, L.; and Yuan, L. 2022a. Bevt: Bert pretraining of video transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 14733–14743.

Wang, R.; Chen, D.; Wu, Z.; Chen, Y.; Dai, X.; Liu, M.; Yuan, L.; and Jiang, Y.-G. 2023b. Masked video distillation: Rethinking masked feature modeling for self-supervised video representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 6312–6322.

Wang, Y.; Li, K.; Li, Y.; He, Y.; Huang, B.; Zhao, Z.; Zhang, H.; Xu, J.; Liu, Y.; Wang, Z.; et al. 2022b. Internvideo: General video foundation models via generative and discriminative learning. *arXiv preprint arXiv:2212.03191*.

Wei, C.; Fan, H.; Xie, S.; Wu, C.-Y.; Yuille, A.; and Feichtenhofer, C. 2022. Masked feature prediction for self-supervised visual pre-training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 14668–14678.

Xie, Z.; Zhang, Z.; Cao, Y.; Lin, Y.; Bao, J.; Yao, Z.; Dai, Q.; and Hu, H. 2022. SimMIM: A Simple Framework for Masked Image Modeling. arXiv:2111.09886.

Yang, J.; Dong, X.; Liu, L.; Zhang, C.; Shen, J.; and Yu, D. 2022. Recurring the transformer for video action recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 14063–14073.

Yang, R.; Mentzer, F.; Van Gool, L.; and Timofte, R. 2021. Learning for Video Compression With Recurrent Auto-Encoder and Recurrent Probability Model. *IEEE Journal of Selected Topics in Signal Processing*, 15(2): 388–401.

Yao, G.; Lei, T.; and Zhong, J. 2019. A review of Convolutional-Neural-Network-based action recognition. *Pattern Recognition Letters*, 118: 14–22. Cooperative and Social Robots: Understanding Human Activities and Intentions.

Ye, H.; Wu, Z.; Zhao, R.-W.; Wang, X.; Jiang, Y.-G.; and Xue, X. 2015. Evaluating Two-Stream CNN for Video Classification. In *Proceedings of the 5th ACM on International Conference on Multimedia Retrieval*, ICMR '15, 435–442. New York, NY, USA: Association for Computing Machinery. ISBN 9781450332743.

Zhao, Z.; and Liang, P. 2006. A highly efficient parallel algorithm for H. 264 video encoder. In *2006 IEEE International Conference on Acoustics Speech and Signal Processing Proceedings*, volume 5, V–V. IEEE.

Zheng, Z.; Yang, L.; Wang, Y.; Zhang, M.; He, L.; Huang, G.; and Li, F. 2023. Dynamic Spatial Focus for Efficient Compressed Video Action Recognition. *IEEE Transactions on Circuits and Systems for Video Technology*.