# $F^3$-Pruning: A Training-<u>F</u>ree and Generalized Pruning Strategy towards <u>F</u>aster and <u>F</u>iner Text-to-Video Synthesis

**Sitong Su**[*], **Jianzhi Liu**[*], **Lianli Gao, Jingkuan Song**

University of Electronic Science and Technology of China (UESTC)
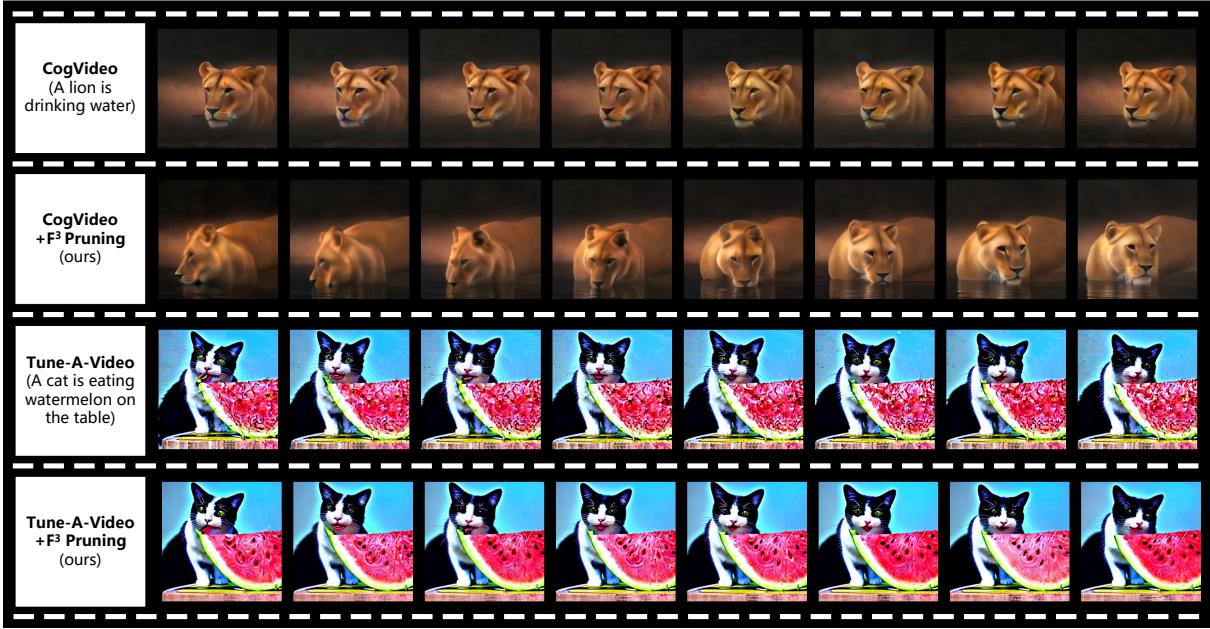sitongsu9796@gmail.com

Figure 1: *Demonstration of some visual results comparison in Text-to-Video (T2V) synthesis*. Our $F^3$-Pruning is applied to the classic transformer-based method CogVideo and the typical diffusion-based method Tune-A-Video. **Without any extra training**, $F^3$-Pruning not only boosts inference efficiency of T2V but also enhances video quality. On the public video dataset UCF-101, applying $F^3$-Pruning to CogVideo makes it **1.35x** faster and promotes video quality metrics FVD by **22%**.

## Abstract

Recently Text-to-Video (T2V) synthesis has undergone a breakthrough by training transformers or diffusion models on large-scale datasets. Nevertheless, inferring such large models incurs huge costs. Previous inference acceleration works either require costly retraining or are model-specific. To address this issue, instead of retraining we explore the inference process of two mainstream T2V models using transformers and diffusion models. The exploration reveals the redundancy in temporal attention modules of both models, which are commonly utilized to establish temporal relations among frames. Consequently, we propose a training-free and generalized pruning strategy called **$F^3$-Pruning** to prune redundant temporal attention weights. Specifically, when aggregate temporal attention values are ranked below a certain ratio, corresponding weights will be pruned.1 Extensive experiments on three datasets using a classic transformer-based model CogVideo and a typical diffusion-based model Tune-

A-Video verify the effectiveness of $F^3$-Pruning in inference acceleration, quality assurance and broad applicability.

## 1 Introductions

Text-to-Video (T2V) synthesis aims to generate high-quality and temporally coherent videos semantically aligned with text inputs. In the early stage of T2V, traditional generative models like GAN (Goodfellow et al. 2014) are trained on datasets of limited scale. Hence, the synthesized videos are restricted by limited semantics and low quality (Li et al. 2018; Pan et al. 2017). To tackle this challenge, recent T2V models utilize large and powerful generative models transformer or diffusion models to train on large-scale datasets (Ho et al. 2022b; Wu et al. 2022a). Alternately, other recent T2V models (Hong et al. 2023; Ho et al. 2022a; Wu et al. 2022b) build upon large-scale T2I models (Ding

et al. 2021; Saharia et al. 2022a; Rombach et al. 2022a) by introducing temporal information, which are also based on transformers or diffusion models. Despite extraordinary progress in video quality and zero-shot generalization ability, inferring such large models incurs huge costs.

Nevertheless, previous inference acceleration methods for transformer or diffusion models either incorporate costly retraining or cannot be generally applied to both two models. Specifically, for transformers-based models, several works are dedicated towards a more efficient transformer (Dosovitskiy et al. 2020; Mehta and Rastegari 2021; Meng et al. 2022). Whereas, these works focus on discriminative tasks without considering synthesis quality. Other works like (Mao et al. 2021; Wei et al. 2023; Rao et al. 2021) prune weights or tokens with costly retraining or finetuning. Differently, ToMe merges tokens without training, which however introduces time-consuming manipulation in inference. For diffusion-based models, traditional diffusion acceleration speeds up the denoising process (Lu et al. 2022), which is limited by the trade-off between video quality and efficiency. Other methods like knowledge distillation (Ho et al. 2022a) or diffusion model pruning (Fang, Ma, and Wang 2023) also incorporate extra training costs. All of the above works are specially designed for one specific model, thus hampering their applicability to different T2V models.

To address the issue, instead of retraining we investigate the inference process of a classic transformer-based model CogVideo (Hong et al. 2023) and a typical diffusion-based model Tune-A-Video (Wu et al. 2022b) to search for the common points. Regardless of model types, the above two models both establish attention between text to each frame, within frames and cross frames to respectively model text-visual alignment, visual quality and temporal coherence as shown in Fig. 2(a). We respectively name these three modules as Cross-modal Attention(CA), Self Attention(SA) and Temporal Attention(TA). Note that, CA, SA and TA are widely adopted in recent T2V models.

Therefore, we statistically analyze the attention values distribution of the three modules during the inference process as shown in Fig. 3. At the bottom of Fig. 3, there are attention maps visualization respectively for CogVideo and Tune-A-Video. The diagonal line of the map refers to SA while the other parts represent TA. As observed, TA is full of values approaching zero, which indicates there exist plenty of non-contributory and redundant parts getting involved in temporal modeling. The redundancy can also be verified by visual results in the first row of Fig. 1 that the motion of the lion is restricted. Moreover, the histograms in Fig. 3 demonstrate the relationship between summed attention values called Aggregate Attention Score(AAS) and network layers or denoising timesteps. As shown, the AAS of TA keeps declining while that of CA or SA increases with the generation process goes. It implies that the importance of generation is gradually transferred from temporal information to visual quality and text-visual alignment. And TA plays a less important role in the late stage of generation. In summary, in the inference stage of both transformers and diffusion models, there exists redundancy in temporal attention modules.

Inspired by the above observation, we propose a training-free and generalized pruning strategy called $F^3$-Pruning to prune redundant and less important temporal attention weights as shown in Fig. 2(b). Instead of designing intrinsic pruning criteria, we claim that attention values could reasonably represent the saliency of corresponding attention weights. Besides, considering the sparsity of TA values, TA values are aggregated into Aggregate Attention Scores (AAS) as pruning criteria rather than complicated comparisons. As demonstrated in Fig. 2(b), our $F^3$-Pruning is applied over network layers or denoising timesteps with the pruning criteria to cut off TA weights whose AAS are ranked below a specific ratio. With $F^3$-Pruning, redundant TA is released and redistributed to SA or CA to further promote video quality and text-visual alignment with coherence.

Our contributions could be summarized as follows:

1) We explore the inference process of two mainstream T2V models using transformer and diffusion, which reveals the redundancy of temporal attention in both models.

2) We propose a training-free and generalized pruning strategy called $F^3$-Pruning to prune redundant temporal attention weights, which both speeds up T2V inference and assures video quality.

3) Extensive experiments on three datasets prove the effectiveness, efficacy and generalization of $F^3$-Pruning. Particularly, $F^3$-Pruning applied to CogVideo on UCF-101 dataset not only speeds up CogVideo by **1.35x** but also significantly improves video quality metric FVD by **22%**.

## 2    Related Works

### 2.1    Text-to-Video Generation

In the initial stage of development, some GAN-based (Li et al. 2018; Pan et al. 2017) models can only create low-resolution videos with restricted semantics. Subsequent works mainly focus on non-adversarial processes. A classic diffusion-based work VDM (Ho et al. 2022b), naturally extends a standard image diffusion model and enables it jointly training with image and video data. Different from that, some works (Wu et al. 2022a; Villegas et al. 2023; Ge et al. 2022) modify a base transformer to adapt to the video synthesis involving multi-model signals.

However, these works require exceptionally high training costs for training from scratch. Therefore some researchers explore the utilization of pretrained T2I models to ease training costs for T2V generation. For instance, CogVideo (Hong et al. 2023) devises a large-scale T2I transformer CogView2 (Ding et al. 2022) by introducing the temporal information through the use of attention among frames. On a different note, Make-A-Video (Singer et al. 2023) breaks away from the usual reliance on text-video pairs for T2V generation by utilizing a pretrained T2I model. Imagen Video (Ho et al. 2022a) follows Imagen (Saharia et al. 2022b), employing a cascaded diffusion model with attention and convolution at multiple resolutions. Moreover, as the quality of video generation improves, recent works begin to explore diverse settings of generation. Tune-A-Video (Wu et al. 2022b) proposes a one-shot video tuning approach for T2V generation, incorporating temporal attention into
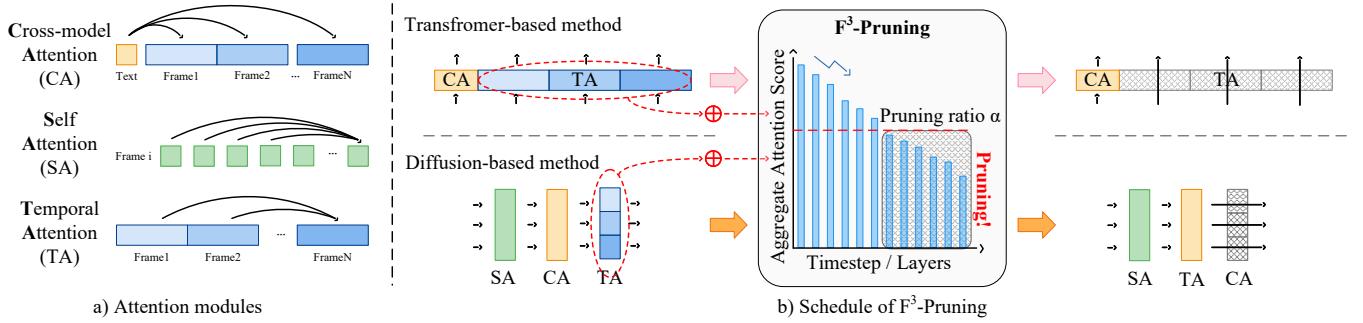
Figure 2: *Overview of our proposed $F^3$-Pruning*. In a), we show three attention modules Cross-model Attention (CA), Self Attention (SA) and Temporal Attention (TA), which are commonly used in T2V to respectively model text-visual alignment, visual quality within each frame and temporal coherence among frames. In b), we demonstrate the schedule of our $F^3$-Pruning applied to the transformer-based methods and the diffusion-based methods. TA weights will be pruned when the sums of TA values of some network layers or denoising timesteps, called Aggregate Attention Score, are ranked below a pruning ratio $\alpha$.

Stable diffusion (Rombach et al. 2022b). And Text2Video-Zero (Khachatryan et al. 2023) enables zero-shot T2V generation without training video.

Since these methods primarily prioritize video quality, semantical consistency, and coherence, they overlook the issue of huge computational costs and long latency in reference.

## 2.2 Inference Acceleration

The methods for model acceleration at the algorithmic level are mainly divided into pruning (Xia, Zhong, and Chen 2022; Lagunas et al. 2021; Zhuo et al. 2022), quantization (Qin et al. 2022; Li et al. 2022; Liu et al. 2022; Song et al. 2019). In this work, we focus on pruning, which is classified into structural and unstructured pruning. Some works focus on unstructured pruning (Dong, Chen, and Pan 2017; Park et al. 2020; Sanh, Wolf, and Rush 2020; Lee et al. 2020), which finds the unimportant parameters and conceals them by masking or setting them to zero. Those methods usually do not bring actual acceleration. Whereas very recent works perform on structural pruning (Ding et al. 2019; You et al. 2019; Liu et al. 2021), which stands out due to its capability to physically remove parameters and substructures from neural networks for accelerating training or inference. Nevertheless, the above research on network pruning has predominantly concentrated on tasks involving discrimination, notably classification endeavors and these approaches all require additional training to determine the locations for pruning.

To avoid training costs, Bolya et al. (Bolya and Hoffman 2023; Bolya et al. 2023) propose a token merging method that notably enhances inference speed by reducing tokens through cosine similarity, all without requiring any training. Nevertheless, this approach entails extra computations to determine which tokens should be merged. This can potentially add substantial overhead to the inference process. Our approach is different from prior methods as we harness temporal redundancy, negating the necessity for intricate token similarity comparisons and sidestepping supplementary computations during inference. In this study, we change the token-specific operations to layer-specific operations by an-

alyzing temporal redundancy and demonstrate the feasibility of this pruning approach in video generation.

## 3 Methodology

Firstly, we introduce some preliminaries in T2V in Sec. 3.1 for further illustration. Then in Sec. 3.2, we investigate the inference stage of two mainstream T2V models to analyze the redundant parts for pruning. Based on the analysis, we introduce our pruning strategy $F^3$-Pruning in Sec. 3.3.

## 3.1 Preliminary

Recent Text-to-Video models are mainly established on powerful generative models transformers and diffusion models. Despite the model types, Cross-modal Attention (CA), Self Attention (SA) and Temporal Attention (TA) are commonly utilized to respectively assure text-visual alignment, visual quality of each single frame and temporal coherence among frames. Specifically, as depicted in Fig. 2 a), CA models attention between text features $T$ and visual features of each frame $F_i$. And SA establishes attention matrices within each frame for visual quality. Learning attention assignation among different frames forms TA for temporal modeling. Consequently, the formulations of these three attention modules could be summarized as follows:

$$\begin{cases} CA = Attn(T, F_i), \\ SA = Attn(F_i, F_i), \quad i,j \in (1,...,N) \\ TA = Attn(F_i, F_j), \end{cases} \quad (1)$$

where $N$ refers to the number of generated video frames. When $i = j$, SA is contained in TA. According to the formulations, the time complexity of TA is $\mathcal{O}(N^2)$ in contrast to $\mathcal{O}(N)$ of SA, which indicates a quadratic time consumption with the increas of frames or resolution.

For the transformer-based methods, we consider CogVideo (Hong et al. 2023), a classic large-scale open-source T2V model as an example. As depicted on the left top of Fig. 2 b), all video frames and text are discretized into tokens and different attentions are entangled to jointly model text and visual connections as
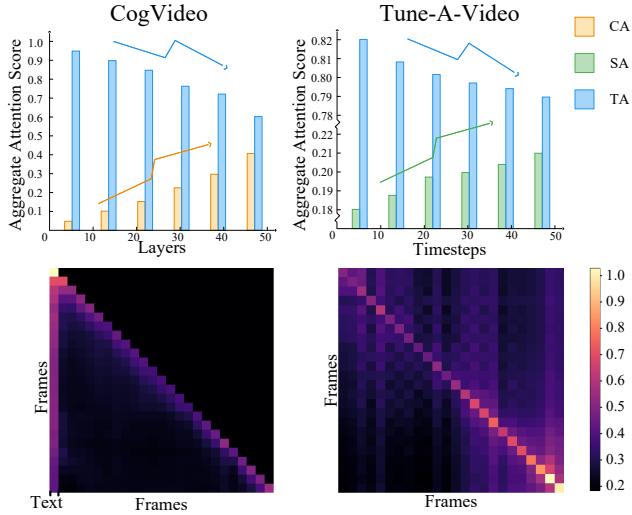
$$[CA; TA], \quad (2)$$

Figure 3: **Top**: Demonstration of the relation between Aggregate Attention Score (AAS) and network layers or denoising timesteps. AAS is declining with the inference step. **Bottom**: Attention Visualization. The diagonal line represents SA, and the upper and lower triangles represent TA. In particular, the leftmost bright line in CogVideo represents CA. As seen, attention values are sparsely distributed.

where SA is contained in TA.

For the diffusion-based methods, we select Tune-A-Video (Wu et al. 2022b) which serves as a baseline one-shot T2V model. As shown on the left bottom of Fig. 2 b), in each denoising step, CA, SA and TA are disentangled and arranged in a cascaded way as

$$TA(CA(SA)). \quad (3)$$

## 3.2 What to Prune: Redundancy Analysis

To avoid costly retraining or finetuning, we probe into the inference stage of the above two representative T2V models CogVideo and Tune-A-Video. Since CA, SA and TA are commonly shared among T2V models, we choose to analyze their attention values distribution in inference, which are summarized into the following three observations.

**Observation 1.** *A large portion of temporal attention values approaching zero indicates redundancy for temporal modeling.* For direct observation, we visualize attention maps of CogVideo and Tune-A-Video in the bottom of Fig. 3. Specifically, for the attention map of CogVideo in the left bottom, values in the most left column, the diagonal line and the other parts in the lower triangular matrix respectively represent CA, SA and TA. As seen, attention values are intensively assigned to CA and SA while sparsely distributed to TA. Moreover, in TA the neighboring frames are more dedicated compared to more previous frames. The conclusion is also compatible with the common sense that a frame is mostly related to its neighboring frames or itself rather than attending to all other frames. In addition, for the attention map visualization of Tune-A-Video in the right bottom of Fig. 3, the diagonal line refers to SA and the other parts

refer to TA. Also, SA is much more emphasized compared to TA, which is consistent with the observation in CogVideo.

Generally, the distribution of attention values indicates that analyzing linguistic information and building visual information within a frame matters much more than attending to all the precedent frames. Based on the observation above, we can conclude that there exists a large extent of redundancy for fully connecting each two frames in TA. The redundancy can also be verified by visual results of CogVideo in Fig. 1. The dense connection among frames restricts the motion of the video, where the lion is kept unchanged instead of drinking water.

**Observation 2.** *The decrease of aggregate temporal attention values with generation process implies its declining importance.* For further analysis, we utilize text prompts from public datasets, which are fed into CogVideo and Tune-A-Video to infer corresponding videos. Then in different network layers or denoising timesteps, we statistically collect the sum of attention values in CA, SA and TA, which are called Aggregate Attention Score (AAS). The corresponding results form the histogram in the top of Fig. 3. As observed, with the deepening network layers in the transformer or increasing denoising timesteps in the diffusion model, the AAS keeps decreasing in TA in contrast with the rise of CA or SA. This implies the decreasing importance of TA along with the generation process. For CogVideo, the increasing AAS in CA demonstrates that generation is more oriented towards textual information, which is consistent with the analysis in CogVideo. As for Tune-A-Video, as the denoising process goes, the focus of generation attends more and more to visual details within a frame rather than temporal relations among frames. The observation also corresponds to a conclusion of diffusion models that the early denoising process focuses on low-frequency contents while the late one focuses more on high-frequency details.

The overall observations reveal that temporal modeling takes different proportions in the whole generation process, and matters less important in the late stage of generation.

**Observation 3.** *Removing redundant TA leads to reasonable attention redistribution to CA or SA, thus contributing to video quality.* For CogVideo, TA and CA are entangled modeling as shown in the left top of Fig. 2(b), which means that their normalized attention values are summed to be 1. Once TA is released, its attention values will be redistributed to CA, thus contributing to text-visual alignment. For Tune-A-Video, CA, SA and TA are cascaded connected as described in Eq. 3. When the overly tight restriction of TA is removed, the influence of SA and CA will be enlarged correspondingly. Especially, if less important TA weights in the late stage of the denoising process are pruned, focusing more on SA will favor the construction of high-frequency details. This observation can be verified through visual results of Tune-A-Video+F³-Pruning as shown in Fig. 1.

In summary, observations 1 and 2 indicate the redundancy of TA and decreased importance of TA in the late stage of generation. Moreover, observation 3 further analyzes the merits of releasing redundant TA, which will contribute to video quality or text-visual alignment.

| Method | UCF-101 | | DAVIS | | LONGTEXT | Performance | |
|---|---|---|---|---|---|---|---|
| | FVD ↓ | Clip(text) ↑ | FVD ↓ | Clip(text) ↑ | Clip(text) ↑ | PFLOPS ↓ | Time (s) ↓ |
| CogVideo | 2537 | **25.03** | 1851 | **26.34** | 28.07 | 9484 | 486 |
| + Pruning (*) | 4073 | 24.66 | 2239 | 25.46 | 27.94 | **5271** | 672 |
| + Reorganization (*) | 4091 | 24.61 | 2037 | 25.81 | 27.73 | 5280 | 688 |
| + TPS (*) | 4471 | 24.35 | 2378 | 25.16 | 27.83 | 5280 | 704 |
| + ToMe (*) | 3594 | 24.74 | 1837 | 26.34 | 28.07 | 5290 | 985 |
| + F$^3$-Pruning | **1990** | 24.77 | **1781** | 26.21 | **28.4** | 5311 | **359** |

Table 1: Quantitative comparison results on five pruning methods applied to the typical transformer-based model CogVideo. "*" represents that the method itself is not originally designed for auto-regressive transformers, and we adapt it to fit CogVideo.

## 3.3 How to Prune: F$^3$-Pruning

Motivated by the observation and analysis from Sec. 3.2, we propose F$^3$-Pruning to structurally prune weights of redundant temporal attention as shown in Fig. 3(b).

**Aggregate Attention Score.** Instead of designing complicated pruning criteria that usually demand extra heads to be trained, we claim that attention values could act as a reasonable indication of the corresponding weights importance. Considering the sparsity of temporal attention, bits-by-bits attention values comparison or ranking will not only cause a large amount of calculations but also hurts the robustness due to the variance of comparing vectors with small norms. In contrast, comparing the sum of temporal attention values avoids dense matrix calculations and could aggregate the sparse values to better represent the importance of the overall temporal attention. Consequently, we propose to sum over the temporal attention values for each network layer or denoising timestep as our pruning criteria coined as Aggregate Attention Score (AAS). The corresponding formulation could be summarized as follows:

$$AAS = sum(Attn(F_i, F_j)) \quad i \neq j. \quad (4)$$

**F$^3$-Pruning.** Given the criteria, our F$^3$-Pruning schedule is illustrated in Fig. 3(b). Specifically, we iterate through all the training samples in a public video dataset and infer corresponding video samples using the given T2V model. Afterward, we calculate over each network layer or denoising timestep following Eq. 4 to obtain the AAS. According to the hypothesis that 'smaller-norm-less-important', we rank the obtained AAS for the whole generation process by their norms. Then we prune the temporal attention weights with 'smaller-norm' below a fixed threshold. Note that different T2V models own different extents in temporal redundancy, which means a fixed threshold will not be adaptable. Therefore, we choose to prune a proportion of temporal attention weights defined by a pruning ratio $\alpha$, where $0 \leq \alpha \leq 1$. That is to say, we prune $\alpha$ of temporal attention weights whose AAS rank in the last $\alpha$ part. Regarding the decreasing importance of temporal attention, we actually prune the TA weights in the late stage of generation.

## 4 Experiments and Analysis

### 4.1 Datasets and Evaluation Metrics

**Datasets** To prove the effect of F$^3$-Pruning, and compare it to other SOTA pruning methods. We follow the set-

| Method | Clip (img) ↑ | Clip (text) ↑ | TFLOPS ↓ | Time (s) ↓ |
|---|---|---|---|---|
| Tune-A-Video | 97.78 | 32.83 | 379.8 | 45.08 |
| + ToMeSD | 97.73 | 32.88 | 369.53 | 48.54 |
| + F$^3$-Pruning | **97.95** | **33.30** | **355.23** | **43.83** |

Table 2: Quantitative comparison results on two pruning methods applied to the typical diffusion-based models Tune-A-Video. Experiment conducts on LONGTEXT.

tings of (Hong et al. 2023; Wu et al. 2022b) and select two public video datasets, including a) **UCF-101**, an action recognition dataset (Peng, Zhao, and Zhang 2019) and b) **DAVIS**, a densely annotated video segmentation (Voigtlaender et al. 2020), which both contain label-to-video pairs. Furthermore, to address the potential loss of text information from using labels as inputs, we build a c) **LONGTEXT** dataset, which consists of full sentences generated by ChatGPT(OpenAI 2021).

**Evaluation Metrics** To comprehensively evaluate the effectiveness and efficacy of different pruning methods, we adopt four commonly used metrics as follows. We utilize a) **FVD** (Peng, Zhao, and Zhang 2019), which is pretained by SVC (StarCraft 2 Videos) dataset, to evaluate the quality of spatial and temporal features in video generation, b) **Clip-Score** (Hessel et al. 2021) to measure the alignment of text and image and coherence of generated videos respectively denoted as Clip(text) and Clip(img), c) **FLOPS** (floating point operations per second) to assess the floating calculation of pruning techniques, and d) **time**(s) to provide an assessment of the practical running speed.

### 4.2 Comparison on transformer-based methods

To verify the efficacy and effectiveness of F$^3$-Pruning, we conduct comparative experiments on UCF-101, DAVIS, and LONGTEXT using the baselines including Token Pruning (Rao et al. 2021), Token Reorganization (Mao et al. 2021), TPS (Wei et al. 2023), and ToMe (Bolya et al. 2023). We equally prune 50% for each baseline on TA. The quantitative and qualitative results are respectively shown in Tab. 1 and Fig. 4.

As indicated in Tab. 1, we equally prune 50% of TA in each baseline and make the FLOPS decline about 44%. However, only F$^3$-Pruning achieves the lowest inference time, resulting in an approximately 1.35x speedup. On the
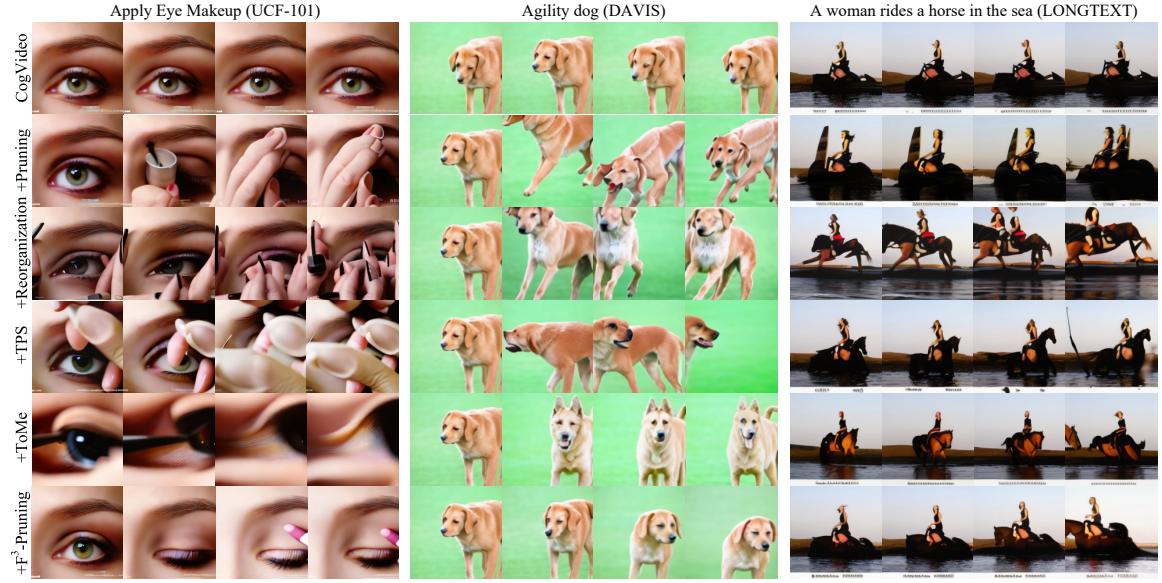
Figure 4: Some examples generated by five pruning methods applied to CogVideo. As demonstrated, $F^3$-Pruning performs the best in coherence and text comprehension.



Figure 5: Some examples generated by two pruning methods applied to the Tune-A-Video (TAV) on the datasets of LONG-TEXT. As demonstrated, $F^3$-Pruning performs the best, especially in object details.

contrary, other baselines lead to an increase in inference time. It is because the $F^3$-Pruning without additional calculations and tensor merging during the inference stage. For video quality, $F^3$-Pruning exhibits exceptional quality on UCF-101, leading to a significant 22% improvement in FVD. Similarly, $F^3$-Pruning achieves the best Clip(text) among the five baselines. Although the Clip(text) of ToMe is slightly higher than $F^3$-Pruning, the results of DAVIS dataset also show the same conclusion as UCF-101. To test the performance of $F^3$-Pruning with natural text as input. We conduct experiments on LONGTEXT dataset and prove that $F^3$-Pruning achieves the best Clip(text) among the five baselines.

The comparison in Fig. 4 shows that four competitors lead to misinterpretations of text and the generation of illogical videos. In the case of the CogVideo, it overly con-

nects dense TA to ensure smooth transitions between frames. However, this design might result in reduced motion information within the generated videos. Additionally, incorrect pruning and merging strategies by the other four competitors lead the misleading text and incoherent video. When $F^3$-Pruning prunes temporal redundancy, the model redistributes TA from previous frames to the textual content. Consequently, this adjustment tends to produce coherent videos with a greater emphasis on text motion alignment.

### 4.3 Comparison on the diffusion-based method

To further verify the versatility of $F^3$-Pruning, we conduct experiments on the diffusion-based Tune-A-Video, with the comparison method ToMeSD (Bolya and Hoffman 2023). As shown in Tab. 2, $F^3$-Pruning achieves the best overall metrics. Note that, although ToMeSD has a little improve-

Figure 6: Some examples of pruning TA, SA, and CA.

| Methods | FVD ↓ | Clip(text) ↑ | PFLOPS ↓ | Time(s) ↓ |
|---------|-------|--------------|----------|-----------|
| CogVideo | 5357 | **25.03** | 9483 | 486 |
| w/o CA | 4730 | 24.97 | 9446 | 480 |
| w/o SA | 5410 | 24.87 | 8651 | 460 |
| w/o TA | **4239** | 24.77 | **5311** | **359** |

Table 3: Abaltion study of w/o CA, SA and TA on UCF-101.



Figure 7: Some examples of different pruning ratios $\alpha$.

| Methods | FVD ↓ | Clip(text) ↑ | PFLOPS ↓ | Time (s) ↓ |
|---------|-------|--------------|----------|------------|
| CogVideo | 5357 | **25.03** | 9483 | 486 |
| $\alpha$=25% | 4937 | 24.96 | 7397 | 422 |
| $\alpha$=50% | 4239 | 24.77 | 5310 | 359 |
| $\alpha$=75% | **1750** | 23.52 | **3224** | **311** |

Table 4: Ablation Study of pruning ratios $\alpha$ on UCF-101.

ment on FLOPS, the inference time of it is slightly higher than Tune-A-Video. It is because that ToMeSD introduces an extra step of merge and unmerge process, which leads to an impact on the efficiency of generation. Additionally, ToMeSD exhibits a slight decrease in Clip(img) and a minor increase in Clip(text), but $F^3$-Pruning demonstrates a more pronounced improvement in Clip(text) and Clip(img). The visual results in Fig. 5 also show that the excessive attention of Tune-A-Video and ToMeSD to other frames causes overly reliance on preceding frames, resulting in distortions and increased noises in the details. As a result, $F^3$-Pruning is the clear winner in this experiment.

## 4.4 Ablation Study

**Effect of pruning CA, SA, and TA** To validate the effects of three attention, we conduct an ablation study on UCF-101 dataset by separately pruning CA, SA, and TA. The results are demonstrated in Tab. 3. Although "w/o CA" has the improvement of FVD and the lowest harmful for Clip(text), its inference time has little improvement. Similarly, "w/o SA" has modest harm for Clip(text), but it increases FVD, thereby reducing the quality of generated video. Differently, "w/o CA" significantly reduced inference time, and making the best trade-off between costs and quality among those three. Our visual results, depicted in Fig. 6, underscore the consequence, revealing "w/o CA" and "w/o SA" have a loss of alignment from textual input. As a result, we choose the
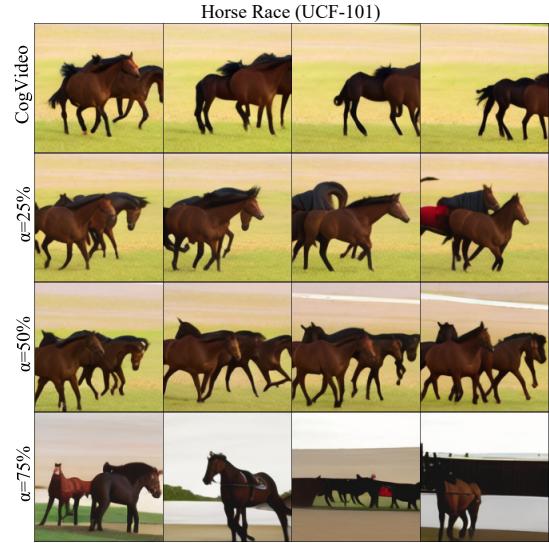
TA as our pruning attention.

**Effect of different pruning ratios** To determine the optimal pruning ratio $\alpha$, we conduct an ablation study to prune CogVideo by 25%, 50%, and 75% on UCF-101. Qualitative results are demonstrated in Tab. 4). Notably, as the pruning ratio increases, we observe a huge reduction in inference time and FLOPS. These there pruning ratios lead to a substantial enhancement in FVD. And 25% and 50% exhibit fewer effects on the final Clip(text), suggesting they maintain a robust text-image alignment. However, the experiment of 75% pruning ratio achieves the best FVD performance and inference time, but it has a notable decline in text-video alignment. This observation is further supported by the findings presented in Fig. 7. As shown in the figure, a 50% pruning ratio demonstrates modest improvements in text-image alignment and inter-frame coherence within the generated videos. In contrast, the experiment with a 75% pruning ratio exhibits a significant loss of coherence among frames. As a result, we select a 50% pruning ratio for all our experiments.

## 5 Conclusion

In this paper, to accelerate inference of T2V, we propose a training-free and generalized pruning strategy called $F^3$-Pruning. Specifically, temporal attention weights will be pruned if their aggregate attention values rank below a pruning ratio. Experiments prove that our method promotes inference speed and video quality on mainstream T2V models.

# References

Bolya, D.; Fu, C.; Dai, X.; Zhang, P.; Feichtenhofer, C.; and Hoffman, J. 2023. Token Merging: Your ViT But Faster. In *ICLR*.

Bolya, D.; and Hoffman, J. 2023. Token Merging for Fast Stable Diffusion. *CoRR*, abs/2303.17604.

Ding, M.; Yang, Z.; Hong, W.; Zheng, W.; Zhou, C.; Yin, D.; Lin, J.; Zou, X.; Shao, Z.; Yang, H.; and Tang, J. 2021. CogView: Mastering Text-to-Image Generation via Transformers. In *NeurIPS*.

Ding, M.; Zheng, W.; Hong, W.; and Tang, J. 2022. CogView2: Faster and Better Text-to-Image Generation via Hierarchical Transformers. In *NeurIPS*.

Ding, X.; Ding, G.; Guo, Y.; and Han, J. 2019. Centripetal SGD for Pruning Very Deep Convolutional Networks With Complicated Structure. In *CVPR*, 4943–4953.

Dong, X.; Chen, S.; and Pan, S. J. 2017. Learning to Prune Deep Neural Networks via Layer-wise Optimal Brain Surgeon. In *NeurIPS*, 4857–4867.

Dosovitskiy, A.; Beyer, L.; Kolesnikov, A.; Weissenborn, D.; Zhai, X.; Unterthiner, T.; Dehghani, M.; Minderer, M.; Heigold, G.; Gelly, S.; et al. 2020. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In *ICLR*.

Fang, G.; Ma, X.; and Wang, X. 2023. Structural Pruning for Diffusion Models. *CoRR*, abs/2305.10924.

Ge, S.; Hayes, T.; Yang, H.; Yin, X.; Pang, G.; Jacobs, D.; Huang, J.; and Parikh, D. 2022. Long Video Generation with Time-Agnostic VQGAN and Time-Sensitive Transformer. In *ECCV*, 102–118.

Goodfellow, I. J.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A. C.; and Bengio, Y. 2014. Generative Adversarial Networks. *CoRR*.

Hessel, J.; Holtzman, A.; Forbes, M.; Bras, R. L.; and Choi, Y. 2021. CLIPScore: A Reference-free Evaluation Metric for Image Captioning. In *EMNLP*, 7514–7528.

Ho, J.; Chan, W.; Saharia, C.; Whang, J.; Gao, R.; Gritsenko, A. A.; Kingma, D. P.; Poole, B.; Norouzi, M.; Fleet, D. J.; and Salimans, T. 2022a. Imagen Video: High Definition Video Generation with Diffusion Models. *CoRR*, abs/2210.02303.

Ho, J.; Salimans, T.; Gritsenko, A. A.; Chan, W.; Norouzi, M.; and Fleet, D. J. 2022b. Video Diffusion Models. In *NeurIPS*.

Hong, W.; Ding, M.; Zheng, W.; Liu, X.; and Tang, J. 2023. CogVideo: Large-scale Pretraining for Text-to-Video Generation via Transformers. In *ICLR*.

Khachatryan, L.; Movsisyan, A.; Tadevosyan, V.; Henschel, R.; Wang, Z.; Navasardyan, S.; and Shi, H. 2023. Text2Video-Zero: Text-to-Image Diffusion Models are Zero-Shot Video Generators. *CoRR*, abs/2303.13439.

Lagunas, F.; Charlaix, E.; Sanh, V.; and Rush, A. M. 2021. Block Pruning For Faster Transformers. In *EMNLP*, 10619–10629.

Lee, N.; Ajanthan, T.; Gould, S.; and Torr, P. H. S. 2020. A Signal Propagation Perspective for Pruning Neural Networks at Initialization. In *ICLR*.

Li, Y.; Min, M. R.; Shen, D.; Carlson, D. E.; and Carin, L. 2018. Video Generation From Text. In *AAAI*, 7065–7072.

Li, Z.; Yang, T.; Wang, P.; and Cheng, J. 2022. Q-ViT: Fully Differentiable Quantization for Vision Transformer. *CoRR*, abs/2201.07703.

Liu, L.; Zhang, S.; Kuang, Z.; Zhou, A.; Xue, J.; Wang, X.; Chen, Y.; Yang, W.; Liao, Q.; and Zhang, W. 2021. Group Fisher Pruning for Practical Network Compression. In *ICML*, 7021–7032.

Liu, Z.; Oguz, B.; Pappu, A.; Xiao, L.; Yih, S.; Li, M.; Krishnamoorthi, R.; and Mehdad, Y. 2022. BiT: Robustly Binarized Multi-distilled Transformer. In *NeurIPS*.

Lu, C.; Zhou, Y.; Bao, F.; Chen, J.; Li, C.; and Zhu, J. 2022. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps. *NeurIPS*.

Mao, J.; Xue, Y.; Niu, M.; Bai, H.; Feng, J.; Liang, X.; Xu, H.; and Xu, C. 2021. Voxel Transformer for 3D Object Detection. In *ICCV*, 3144–3153.

Mehta, S.; and Rastegari, M. 2021. MobileViT: Lightweight, General-purpose, and Mobile-friendly Vision Transformer. In *ICLR*.

Meng, L.; Li, H.; Chen, B.-C.; Lan, S.; Wu, Z.; Jiang, Y.-G.; and Lim, S.-N. 2022. Adavit: Adaptive vision transformers for efficient image recognition. In *CVPR*, 12309–12318.

OpenAI. 2021. ChatGPT. https://openai.com/research/chatgpt.

Pan, Y.; Qiu, Z.; Yao, T.; Li, H.; and Mei, T. 2017. To Create What You Tell: Generating Videos from Captions. In *ACM MM*, 1789–1798.

Park, S.; Lee, J.; Mo, S.; and Shin, J. 2020. Lookahead: A Far-sighted Alternative of Magnitude-based Pruning. In *ICLR*.

Peng, Y.; Zhao, Y.; and Zhang, J. 2019. Two-Stream Collaborative Learning With Spatial-Temporal Attention for Video Classification. *TCSVT*, 29(3): 773–786.

Qin, H.; Ding, Y.; Zhang, M.; Yan, Q.; Liu, A.; Dang, Q.; Liu, Z.; and Liu, X. 2022. BiBERT: Accurate Fully Binarized BERT. In *ICLR*.

Rao, Y.; Zhao, W.; Liu, B.; Lu, J.; Zhou, J.; and Hsieh, C. 2021. DynamicViT: Efficient Vision Transformers with Dynamic Token Sparsification. In *NeurIPS*, 13937–13949.

Rombach, R.; Blattmann, A.; Lorenz, D.; Esser, P.; and Ommer, B. 2022a. High-resolution image synthesis with latent diffusion models. In *CVPR*.

Rombach, R.; Blattmann, A.; Lorenz, D.; Esser, P.; and Ommer, B. 2022b. High-Resolution Image Synthesis with Latent Diffusion Models. In *CVPR*, 10674–10685.

Saharia, C.; Chan, W.; Saxena, S.; Li, L.; Whang, J.; Denton, E. L.; Ghasemipour, K.; Gontijo Lopes, R.; Karagol Ayan, B.; Salimans, T.; et al. 2022a. Photorealistic text-to-image diffusion models with deep language understanding. *NeurIPS*, 35: 36479–36494.

Saharia, C.; Chan, W.; Saxena, S.; Li, L.; Whang, J.; Denton, E. L.; Ghasemipour, S. K. S.; Lopes, R. G.; Ayan, B. K.; Salimans, T.; Ho, J.; Fleet, D. J.; and Norouzi, M. 2022b. Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding. In *NeurIPS*.

Sanh, V.; Wolf, T.; and Rush, A. M. 2020. Movement Pruning: Adaptive Sparsity by Fine-Tuning. In *NeurIPS*.

Singer, U.; Polyak, A.; Hayes, T.; Yin, X.; An, J.; Zhang, S.; Hu, Q.; Yang, H.; Ashual, O.; Gafni, O.; Parikh, D.; Gupta, S.; and Taigman, Y. 2023. Make-A-Video: Text-to-Video Generation without Text-Video Data. In *ICLR*.

Song, J.; Zhu, X.; Gao, L.; Xu, X.; Liu, W.; and Shen, H. T. 2019. Deep Recurrent Quantization for Generating Sequential Binary Codes. In *IJCAI*, 912–918.

Villegas, R.; Babaeizadeh, M.; Kindermans, P.; Moraldo, H.; Zhang, H.; Saffar, M. T.; Castro, S.; Kunze, J.; and Erhan, D. 2023. Phenaki: Variable Length Video Generation from Open Domain Textual Descriptions. In *ICLR*.

Voigtlaender, P.; Luiten, J.; Torr, P. H. S.; and Leibe, B. 2020. Siam R-CNN: Visual Tracking by Re-Detection. In *CVPR*, 6577–6587.

Wei, S.; Ye, T.; Zhang, S.; Tang, Y.; and Liang, J. 2023. Joint Token Pruning and Squeezing Towards More Aggressive Compression of Vision Transformers. *CVPR*.

Wu, C.; Liang, J.; Ji, L.; Yang, F.; Fang, Y.; Jiang, D.; and Duan, N. 2022a. NÜWA: Visual Synthesis Pre-training for Neural visUal World creAtion. In *ECCV*, 720–736.

Wu, J. Z.; Ge, Y.; Wang, X.; Lei, W.; Gu, Y.; Hsu, W.; Shan, Y.; Qie, X.; and Shou, M. Z. 2022b. Tune-A-Video: One-Shot Tuning of Image Diffusion Models for Text-to-Video Generation. *CoRR*, abs/2212.11565.

Xia, M.; Zhong, Z.; and Chen, D. 2022. Structured Pruning Learns Compact and Accurate Models. In *ACL*, 1513–1528.

You, Z.; Yan, K.; Ye, J.; Ma, M.; and Wang, P. 2019. Gate Decorator: Global Filter Pruning Method for Accelerating Deep Convolutional Neural Networks. In *NeurIPS*, 2130–2141.

Zhuo, L.; Wang, G.; Li, S.; Wu, W.; and Liu, Z. 2022. Fast-Vid2Vid: Spatial-Temporal Compression for Video-to-Video Synthesis. In *ECCV*.

## Appendix A    Implementation Details

There are two stages for CogVideo. The first stage focuses on the generation of key frames. The second stage focuses on interpolation and super-resolution of the key frames. Since $F^3$-Pruning primarily addresses temporal redundancy during video generation, we specifically select the first stage of CogVideo for the ablation study.

Some details of datasets are described in Tab. 5. All experiments are conducted on NVIDIA A6000 GPUs.

## Appendix B    More Visual Results

To provide additional support for the effectiveness of $F^3$-Pruning, we respectively show more comparison visual results in Fig. 8 and Fig. 9 on UCF-101 and LONGTEXT using CogVideo. Both results show that four competitors generate incoherent and less text-aligned videos. For example, they generate distorted faces and different people in the adjacent frames. And $F^3$-Pruning avoids the above problems to a certain extent and shows the robust text-visual alignment and video coherence as the original CogVideo. And more visual results of $F^3$-Pruning applied to CogVideo on LONGTEXT are shown in Fig. 10.

| Datasets | Text inputs | Batch size | Total videos |
|----------|-------------|------------|--------------|
| UCF-101 | 101 | 10 | 1010 |
| DAVIS | 60 | 8 | 480 |
| LONGTEXT | 50 | 8 | 400 |

Table 5: Details of datasets.

Figure 8: More visual results of different pruning methods on the dataset of UCF-101.



Figure 9: More visual results of different pruning methods on the dataset of LONGTEXT.

A musician playing the harmonica and creating soulful melodies



A cyclist crossing the finish line with arms raised in victory



A person kayaking down a fast-flowing river



A group of friends hiking to the top of a majestic mountain



A cyclist participating in a challenging mountain biking competition



A cyclist racing down a winding country road



A person meditating in a tranquil forest



A woman practicing archery and hitting the bullseye



A group of friends playing beach volleyball in the sand



A child blowing bubbles and chasing them with delight



A musician playing the drums in a lively rock band



A musician playing the accordion and filling the air with lively music



A surfer paddling out into the ocean in search of the perfect wave



A painter creating a colorful mural on a city street



A child riding a roller coaster with excitement and thrill



A group of friends going on a thrilling white-water rafting adventure



A musician playing a soothing melody on a grand piano



A woman writing in her journal while sitting under a shady tree



A woman performing a graceful figure skating routine on ice



A person practicing skateboarding tricks at a skate park



Figure 10: More visual results of F$^3$-Pruning applied to CogVideo on the dataset of LONGTEXT.