RainFusion: Adaptive Video Generation Acceleration via Multi-Dimensional Visual Redundancy

Aiyue Chen^{1*}, Bin Dong^{1*}, Jingru Li¹ Jing Lin¹, Yiwu Yao¹, Gongyi Wang¹ ¹Huawei Technologies Co., Ltd

Abstract

Video generation using diffusion models is highly computationally intensive, with 3D attention in Diffusion Transformer (DiT) models accounting for over 80% of the total computational resources. In this work, we introduce Rain-**Fusion**, a novel training-free sparse attention method that exploits inherent sparsity nature in visual data to accelerate attention computation while preserving video quality. Specifically, we identify three unique sparse patterns in video generation attention calculations-Spatial Pattern, Temporal Pattern and Textural Pattern. The sparse pattern for each attention head is determined online with negligible overhead (~0.2%) with our proposed **ARM** (Adaptive Recognition Module) during inference. Our proposed RainFusion is a plug-and-play method, that can be seamlessly integrated into state-of-the-art 3D-attention video generation models without additional training or calibration. We evaluate our method on leading open-sourced models including HunyuanVideo, OpenSoraPlan-1.2 and CogVideoX-5B, demonstrating its broad applicability and effectiveness. Experimental results show that RainFusion achieves over 2× speedup in attention computation while maintaining video quality, with only a minimal impact on VBench scores (-0.2%).

1. Introduction

Diffusion models have become the leading approach in video generation, demonstrating exceptional performance and broad applicability [1] [9] [28] [24]. Initially built on U-Net architectures [1] [9], the field has transitioned to Diffusion Transformers (DiTs), which now serve as the mainstream approach owing to their enhanced performance and scalability. This architectural evolution has further advanced with the adoption of 3D full-sequence attention mechanisms [24] [28], replacing the previously dominant

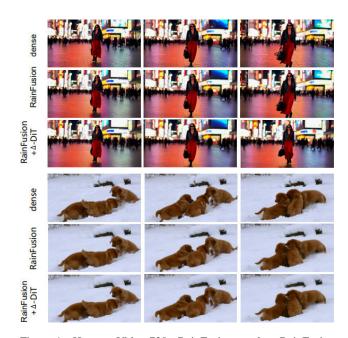


Figure 1. HunyuanVideo 720p RainFusion results. RainFusion and RainFusion combined with Δ -DiT shows good visual quality and high similarity to dense results. Upper prompt: "A stylish woman walks down a Tokyo street filled with warm glowing neon and animated city signage.". Lower prompt: "A litter of golden retriever puppies playing in the snow. Their heads pop out of the snow, covered in.".

2D+1D spatial-temporal attention (STDiT) [21] that separately computes spatial and temporal attention alternatively. Although these advancements have enhanced modeling capabilities, they also impose substantial computational challenges, particularly in attention computation.

The computational complexity of these models scales quadratically with the sequence length, expressed as $O(s^2t^2)$, where s and t represent the spatial and temporal dimensions, respectively. This scaling poses a substantial bottleneck, as evidenced by the deployment of Open-Sora-Plan 1.2 [26] on a single A100 GPU, which requires approximately 48 minutes to generate a 4-second 720p video.

^{*}These authors contributed equally.

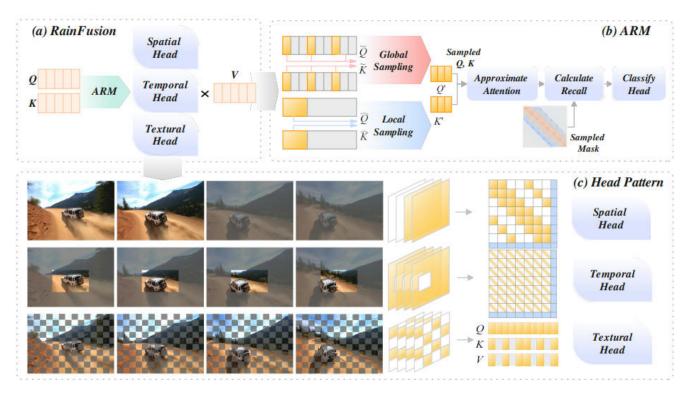


Figure 2. (a) RainFusion pipeline including Adaptive Recognition Module(ARM) and applying sparse pattern to Flash Attention. (b) ARM determine the pattern using subset of query and key to calculate approximates attention score and applying the predefined pattern mask to get attention recall to determine the head category. The sampled queries Q' and keys K' are either sourced from the tokens of the first frame or obtained by sampling from the full set of tokens with equal intervals. (c) The three head sparse pattern. These three heads respectively concentrate on portraying global spatial details with local temporal information, local spatial details with global temporal information, and high-level textural information.

Profiling analysis demonstrates that the attention mechanism consumes over 80% of total computation, making it the principal performance bottleneck in the video generation pipeline.

To improve the computational efficiency of video generation models, researchers have developed two key algorithms: (1) sampling optimization techniques that reduce the number of required inference steps through adaptive sampling schedules [15] [30], and (2) caching mechanisms that exploit redundancy by reusing features across adjacent timesteps [6] [13] [16]. Sampling optimization techniques are inherently limited by their dependency on post-training adjustments, limiting their practical applicability. Furthermore, both sampling optimization and caching algorithms necessitate models to operate with a relatively large number of inference steps, as their effectiveness relies heavily on sufficient redundancy between consecutive timesteps. Despite these advancements, optimizing the attention mechanism has not yet been explored in depth. DiTFastAttn[40] use brute-force sliding-window mask and use residual cache to compensate quality loss. SVG[38] ignores model generality and the inherent visual feature in video.

In this work, we introduce a novel sparse attention mechanism that effectively leverages two key characteristics of video generation: (1) the inherent spatial-temporal redundancy in video, and (2) the importance of specific image texture. We observe that there exists three types of sparse pattern in attention, one for temporal pattern which attends to the same spatial location in different frames, one for spatial pattern which models all spatial location in consecutive frames, the other for detailed texture of video frames. As shown in Fig.3, it is the attention score-map of some heads with the vertical axis and the horizontal axis representing q and k respectively. The first row captures local repetitive patterns within each window, which we define as the Temporal Head, indicating that certain heads consistently attend to the same locations across different frames. The second row reveals more global patterns across neighboring frames, which we term the Spatial Head. The third row highlights **Textural Heads**, where important tokens are attended to by all query tokens. We determine the sparse pattern for each head online using Adaptive Recognition Module(ARM) which only introduce $\frac{1}{t^2}$ overhead where t represents frame number in latent space. The overall pipeline

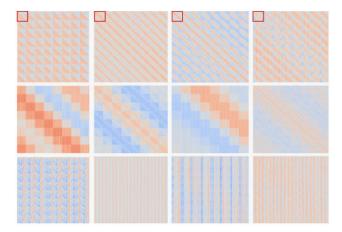


Figure 3. The attention sparsity pattern with the vertical axis and the horizontal axis representing query and key respectively. The first row depicts the temporal sparsity pattern, which models the same spatial location across different frames (with the red box in the upper-left corner highlighting the basic repeated pattern). The second row shows the spatial sparsity pattern, focusing on all locations in neighboring frames. The third row presents a conventional full-attention head, for which we propose a sophisticated textural sparse attention mechanism.

is shown in Fig.2. We first determine sparse pattern of different head online using global or local sampling, and then calculate attention using their respective sparse pattern.

Extensive experiments on different video generation models including OpenSoraPlan-1.2 [26], HunyuanVideo-13B [14], CogVideoX-5B [39] prove the generality and effectiveness of RainFusion. The contributions of this paper include:

- We present RainFusion, a novel plug-and-play framework that leverages tri-dimensional sparsity across spatial, temporal, and textural domains to optimize video diffusion models. The proposed method dynamically determines sparse patterns through online estimation, effectively exploiting the intrinsic redundancy inherent in video data. The name RainFusion is derived from the observation that the sparse patterns resemble the continuous, interconnected lines formed by rain.
- We put forward a simple but potent sparse pattern estimation method ARM that entails minimal computational cost (~0.2% overhead), thereby rendering our RainFusion highly efficient.
- RainFusion can be applied to many SOTA video generation models, OpenSoraPlan-1.2 [26], HunyuanVideo-13B [14], CogVideoX-5B [39] with over 2x speedup in attention at negligible quality loss (-0.2% VBench score) as shown in Fig. 5.

2. Related Work

2.1. Diffusion Models

Diffusion models [4, 8, 10, 18, 22, 25, 42] have surpassed Generative Adversarial Networks (GANs) in generative tasks by iteratively reversing a noisy process to synthesize data, such as images, through progressive denoising. These models typically use U-Net [2, 27, 29] or transformer-based architectures [25], with the latter gaining prominence in vision applications, as seen in DiT (Diffusion Transformers) [25] for data distribution modeling and PixArt- Σ [4] for 4K image generation. Furthermore, diffusion models have been extended to video synthesis [2], with two main approaches emerging: (1) the 2D+1D STDiT structure, used in Open-Sora [42], and (2) the 3D full-sequence attention mechanism, employed by Sora [24], Open-Sora-Plan 1.2 [26], CogVideoX [39], and Hunyuan Video [14]. These developments underscore the versatility and scalability of diffusion models in tackling increasingly complex generative tasks.

2.2. Sparse Attention in Transformers

In Transformer-based large models, the quadratic complexity of matrix multiplication QK^T in attention mechanisms drives high computational costs. To address this, recent research exploits sparsity in attention maps [12, 32, 41], and some research use techniques like token pruning [31, 33, 35] and token merging [3, 34, 37] to reduce sequence length and improve inference efficiency. Some methods employ dynamic sparse attention [12] or merge sparse tokens [32] to accelerate LLM inference. Similarly, in vision-specific models like ViTs and DiTs, sparsity is leveraged through dynamic activation pruning [7], pixel downsampling [31], and KV matrix downsampling [33], while video generation adapts token merging via bipartite soft matching [3], importance sampling [37], and spectrum-preserving techniques [34]. These advancements highlight the broad potential of sparse attention to enhance efficiency across diverse domains.

2.3. Attention Sharing and Cache

When accelerating diffusion model inference, cache methods leverage attention map similarity between adjacent denoising timesteps [16, 19, 36]. For example, Δ-DiT [5] introduces a tailored caching method for DiT acceleration, while DeepCache [20] and TGATE [17] reduce redundant calculations by layer-wise attention similarities. Recent methods further optimize performance by caching model outputs [16] or dynamically adjusting caching strategies [13]. Additionally, techniques like DiTFastAttn [40] combine sparse attention with caching, exploiting spatial, temporal, and conditional redundancies for efficient attention compression. These advancements demonstrate the potential of integrating sparse attention and caching to enhance

the scalability and speed of diffusion model inference.

Recent Work. SVG [38] advances sparse attention research by analyzing spatial and temporal attention sparsity in DiTs and proposing a training-free online profiling strategy. However, they classify attention heads only into temporal and spatial groups, neglecting irregular attention patterns in video generation. Unlike SVG, our work focuses on irregular attention heads to capture fine-grained textural details for improved video generation.

3. Methodology

In this section, we introduce **RainFusion**, a training-free adaptive algorithm, designed to exploit the computing sparsity in 3D full attention to accelerate video generation.

3.1. Preliminary

Existing video generation models utilize 3D full attention mechanisms, which jointly capture both spatial and temporal dependencies to elevate generation quality. However, it comes at a high computational cost.

We define the shape of the latent video as (H,W,T). In 3D full attention, the video sequence is formed by flattening T sub-sequences, each sub-sequence represents a single frame of length $H\times W$. We denote $(Q,K,V)\in\mathbb{R}^{N\times d}$ as the query, key, and value tokens, respectively, and define M as the attention mask with shape $N\times N$, where $N=H\times W\times T$ and d is the hidden dimension of each head. The bidirectional 3D full attention can be formulated as follows:

$$S(Q, K, M) \leftarrow Softmax(\frac{QK^T}{\sqrt{d}} + M)$$
 (1)

$$Attn(Q, K, V, M) \leftarrow S(Q, K, M)V \tag{2}$$

The computation complexity is $O(N^2)$. While 3D full attention mechanisms are inherently dense, our analysis reveals discernible computational sparsity patterns across attention heads. As shown in Fig.3 and Fig. 2 (c), we classify these specialized heads into three categories: Spatial head, Temporal head, and Textural head.

3.2. Attention Head Mechanism Design

Spatial Head The Spatial Head exhibits global spatial dependencies within individual frames while capturing localized temporal dependencies across the full sequence. This characteristic indicates that the Spatial Head emphasizes both the completeness of individual frames and the overall coherence among adjacent or key frames. Consequently, it suggests that certain non-key frames hold relatively less significance and can be excluded from the attention calculation.

$$Attn_{spatial} \leftarrow Attn(Q_f, K_{\{f'\}}, V_{\{f'\}}, M_{spatial})$$
 (3)

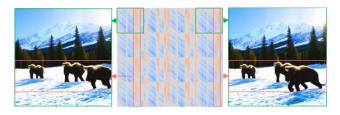


Figure 4. The above figure shows the attention score map of a typical textural head. The green region represents a single frame. Notably, the pink region, characterized by high attention score for most Q, coincides with the motion regions emphasized by the prompt.

Here, $\{f'\}$ denotes the set of significant frames for the f^{th} frame. Therefore, a global striped attention mask $M_{spatial}$ is designed as depicted in Fig.2 (c). A continuous subsequence of a frame is defined as a window segment. The positions of these window segments determine both the key frame attended by the attention mechanism and the resulting computational gains.

Temporal Head Contrary to the Spatial Head, the Temporal Head demonstrates locality within a single-frame subsequence in spatial domain, while exhibits a global characteristic in whole temporal domain. The Temporal Head is particularly attentive to the correlation between the same local regions across different video frames. Its primary focus is on creating regional details that maintain spatial continuity. This unique property can lead to the manifestation of local sparsity within a single-frame sub-sequence and periodic sparsity throughout the entire sequence, as shown in Fig.2 (c).

Textural Head It becomes evident that certain content holds significant importance throughout the entire video, particularly those parts intricately linked to high-level textural description, which is shown in Fig.4. This is manifested in the fact that some specific K, V consistently receive high attention scores for most Q. As a result, while the distribution of tokens is sparse, it is challenging to identify a regular attention mask that can effectively adapt to this sparsity state. Based on the above considerations, we condense the K, V sequence approximately by referring to the property of image downsampling. The K, V sequence will be rearranged and tokens will be retained in a checkerboard interleaving pattern in spatial domain, as depicted in Fig.2 (c).

$$C = \{a_{ij} \mid ((i \bmod \tau = k) \land (j \bmod \tau = k)), \\ 0 \le i < H, 1 \le j < W, 0 \le k < \tau\} \quad (4)$$

$$Attn_{Irregular} \rightarrow Attn(Q, K_{\{C\}}, V_{\{C\}}, M_{init})$$
 (5)

C indicates the set of chosen K,V token indexes, τ represents the stride of the checkerboard, and $M_{\rm init}$ be the allzero mask. The checkerboard format ensures that information from each discarded token in the spatial domain can be implicitly generated by referring to the four nearest remaining tokens. Additionally, we opt to directly retain or discard some tokens rather than averaging them. This is because averaging would obscure the intrinsic information of tokens, making it challenging to implicitly supply the correct information for discarded tokens.

3.3. Adaptive Recognition Module(ARM)

As described in Sec 3.2, RainFusion categorizes all heads into three distinct types: Spatial Head, Temporal Head, and Textural Head. However, we find that the pattern of each head is highly dynamic. For instance, factors such as input prompts and sampling steps all influence the characteristics of each head. Given these considerations, we introduce an Adaptive Recognition Module (ARM). This module is designed to do online and adaptive classification for all heads with minimal computational cost.

We first acquire the approximate attention score, and then compute the masked attention recall. As illustrated in Fig 2 (b), in local sampling, we select the tokens of the first frame sub-sequence as \widehat{Q}, \widehat{K} . In the case of global sampling, we sample tokens at equal intervals ω as \widehat{Q}, \widehat{K} . We utilize the downsampled sequences to calculate the attention score, which serves as an approximation of the overall attention score. Then we compute the masked recall based on the approximate score:

$$R' \leftarrow Recall(Q', K', M') = \frac{S(Q', K', M')}{S(Q', K', M_{init})}$$
 (6)

Q',K' represent the downsampled sequences.M' denotes the attention mask derived by downsampling either $M_{spatial}$ or $M_{temporal}$ in accordance with the corresponding token downsampling rules. S means softmax operation as shown in Equation Attention recall means the proportion of valid information that can be preserved under the current pattern mask. Through this method, we are able to adaptively and efficiently determine the category of each head online with minimal computational overhead.

Algorithm 1 provides a detailed introduction to the process of the Adaptive Recognition Module (ARM).

4. Experiments

4.1. Settings

Models We evaluate RainFusion on three widely adopted video generation models: OpenSoraPlan-1.2 [26], HunyuanVideo-13B [14] and CogVideoX-5B [39]. For HunyuanVideo and OpenSoraPlan-1.2, we generate 125 and 93 frames at 480p resolution, with latent dimensions of

Algorithm 1 Adaptive Recognition Module(ARM)

```
Input Q, K, M_{spatial}, M_{temporal}
Output H # Head Category
\widehat{Q}, \widehat{K}, \widehat{M}_{temporal} \leftarrow LocalSampling(Q, K, M_{temporal})
\widehat{Q}, \widetilde{K}, \widehat{M}_{spatial} \leftarrow GlobalSampling(Q, K, M_{spatial})
\widehat{R} \leftarrow Recall(\widehat{Q}, \widehat{K}, \widehat{M}_{temporal})
\widetilde{R} \leftarrow Recall(\widetilde{Q}, \widetilde{K}, \widehat{M}_{spatial})
if (\widehat{R} \geq \alpha) then
H \leftarrow Temporal Head # high priority for Temporal else if (\widetilde{R} \geq \alpha) then
H \leftarrow Spatial Head
else
H \leftarrow Textural Head
end if
return H
```

(32, 30, 40) and (24, 30, 40) after VAE downsampling and patch embedding, respectively. For CogVideoX-5B, 45 frames are generated at 480×720 , corresponding to a latent shape of (12, 30, 45).

Datasets and Benchmarks VBench [11] is a comprehensive benchmark suite for video generation tasks, systematically decomposing generation quality into 16 distinct evaluation dimensions. It further computes three weighted aggregated scores derived from these dimensions to holistically assess model performance. It consists of 946 prompts for all dimension evaluation. Video generation is a computation-heavy tasks, generating a four second 480p video costs for about 3 minutes. So we only use one random seed instead of five in all the following experiments. In the ablation study, for the sake of accelerating the experiments, we utilize 48 Sora prompts [23]. And we use all VBench 946 prompts when comparing with other methods in section 4.3.

Baselines We show the effectiveness and efficiency of RainFusion to compare it with other sparse or cache-based methods, including DiTFastAttn [40], Δ -DiT[6]. For DiTFastAttn, we use their official configurations. For Δ -DiT, we use the similar accelerate rate of RainFusion. For RainFusion, we set the sparsity to 50% and we keep the first 10% timesteps using dense calculation, which corresponds to about $1.85\times$ speedup in attention. Specifically, we set bandwidth = $\frac{1}{4}$ in both local and global pattern, corresponding to $\frac{9}{16}$ computation reduction. As for textural pattern, we reduce the key value tokens by half using the checkerboard layout as in Section 3.2.

Model	Method	Loss	Quality Score↑	Semantic Score↑	Total Score↑	Speedup
	baseline	/	82.04	70.7	79.77	1.0×
CogvideoX-5B	$\Delta ext{-DiT}$	-5.37	76.56	65.76	74.4	$1.81 \times$
	DiTFastAttn	-5.36	77.94	60.27	74.41	$1.52 \times$
	RainFusion	-0.28	81.64	70.91	79.49	1.85×
OpenSoraPlan-1.2	baseline	/	79.6	38.01	71.28	1.0×
	$\Delta ext{-DiT}$	-0.93	78.51	37.7	70.35	$1.81 \times$
	DiTFastAttn	-2.31	77.56	34.59	68.97	$1.42 \times$
	RainFusion	-0.32	79.08	38.44	70.95	1.91 ×
HunyuanVideo	baseline	/	84.11	70.12	81.31	1.0×
	$\Delta ext{-DiT}$	-0.87	83.27	69.08	80.44	$1.81 \times$
	DiTFastAttn	-1.14	82.87	69.37	80.17	$1.78 \times$
	RainFusion	-0.4	83.77	69.46	80.91	$1.89 \times$
	RainFusion+	-0.19	83.79	70.46	81.12	$1.84 \times$
	RainFusion+ & Δ -DiT	-0.49	83.43	70.35	80.82	$2.37 \times$

Table 1. Comparison with state-of-the-art algorithms.

Model	S	Т	Те	L	Average Loss	Subject Consistency†	Background Consistency†	Motion Smoothness↑	Dynamic Degree↑	Aesthetic Quality↑	Imaging Quality↑
					/	93.48	95.24	97.19	45.83	58.12	64.79
CogvideoX-5B	\checkmark	\checkmark		\checkmark	-1.05	92.87	94.65	97.47	45.83	56.58	60.91
	\checkmark	\checkmark	\checkmark	\checkmark	-0.18	93.27	95.31	97.23	45.83	58.11	63.80
	\checkmark	\checkmark	\checkmark		-0.42	93.08	95.40	97.27	45.83	57.17	63.36
					/	94.65	95.19	99.40	41.67	56.84	57.67
OpenSoraPlan-1.2	\checkmark	\checkmark		\checkmark	-1.29	92.53	94.72	98.94	43.75	55.05	52.66
	\checkmark	\checkmark	\checkmark	\checkmark	1.03	93.22	94.31	99.15	52.08	55.67	57.17
	\checkmark	\checkmark	\checkmark		0.33	92.73	94.69	99.18	50.00	55.12	55.71

Table 2. Ablation Results. RainFusion with three pattern and local estimation achieves the best result. We denote S, T, Te, L as spatial head, temporal head, textural head and use local sampling in estimating local pattern recall, respectively.

4.2. Ablation Study

Component-wise Analysis For RainFusion, there exists three kind of heads (Spatial, Temporal, Textural) as shown in Fig.2 (c). We use ARM in Fig.2 (b) to determine the pattern for each head. We do ablation study on the effectiveness of each head and how to estimate the patterns. We observe that there exists local sparse pattern inner the global sparse pattern as shown in the second row of Fig.3. So if a head is above recall rate for both local and global pattern, we use local pattern first to cover more active region. We use global sampling as described in 3.3 to get global pattern recall. To estimate local pattern recall, we compare the global sampling method and local sampling method(considering only the first frame tokens). As shown in Tab.2, we can see that with all three mask and with local estimation get the best results.

Specifically, for CogVideoX-5B, using only two mask

(same as SVG[38] which only uses spatial and temporal head), the loss is 1.05%. When adding textural head, we can get the best VBench score with the average loss 0.18%. But if we change local estimate to global estimate, the loss is 0.42%. And as shown in Fig.6, we can see that Rain-Fusion using three masks outperform using only two mask similar to SVG[38]. Videos using our method is preserve more details with better imaging quality.

For OpenSoraPlan-1.2, it shows similar results. Using only two mask method drops by 1.29%, when adding textural head, the result even perform better than baseline method. But when using global sampling, the result is a little worse, 0.33% compared to 1.03% improvement with local sampling in VBench score. So we use all three mask and local sampling as our default RainFusion configuration in the following sections.



Figure 5. Video Comparison using CogVideoX-5B with different accelerating algorithms. Left prompt: "A steam train moving on a mountainside." Right Prompt: "a zebra on the left of a giraffe, front view.".



Figure 6. RainFusion video comparisons on CogVideoX-5B. Two Head means only use spatial and temporal head similar to SVG. We can see that RainFusion performs better than SVG and similar to baseline with 1.85x speedup. Left prompt: "Animated scene features a close-up of a short fluffy monster kneeling beside a melting red candle". Right prompt: "A petri dish with a bamboo forest growing within it that has tiny red pandas running around."

Parameter Sensitivity We test different sparse ratio by using different bandwidth and stride in spatial-temporal head and textural head, respectively. We test different Rain-Fusion configuration in CogVideoX-5B of speedup $2.5\times$ and $3.0\times$ by setting the bandwidth of spatial and temporal head to 0.18 and 0.13, and textural stride to 3 and 4, respectively. As shown in Tab.3, our default setting can achieve $1.85\times$ speedup with 0.21% average loss, while for $2.5\times$ and $3.0\times$ speedup, the loss is 0.56% and 1.35%. The $2.5\times$ RainFusion performs better than only local and global mask shown in Tab.2 first row of CogVideoX-5B part, sim-

ilar to SVG(1.85 \times speedup). As for 3.0 \times speedup, the loss is 1.35%. It shows that our method can achieve more speedup with a tradeoff of accuracy. We can conclude that our method shows robust performance for different speedup ratio as shown in Fig.7.

4.3. Comparison with Baselines

Quantitative Results We compare RainFusion with sparse method DiTFastAttn and cache-based method Δ -DiT. The result is shown in Table 1. We can see that in similar speedup, RainFusion performs best among other accel-

Method	Average Loss	Subject Consistency†	Background Consistency [↑]	Motion Smoothness†	Aesthetic Quality↑	Imaging Quality↑
baseline	/	93.48	95.24	97.19	58.12	64.79
1.85× RainFusion	-0.21	93.27	95.31	97.23	58.11	63.81
2.50× RainFusion	-0.56	92.80	95.08	97.14	57.56	63.40
3.00× RainFusion	-1.36	92.36	94.85	97.27	55.20	62.30

Table 3. Parameter Sensitivity Experiment Results.



Figure 7. Video Comparison using CogVideoX-5B with different speedup ratio. Left prompt: "A cat waking up its sleeping owner demanding breakfast." Right Prompt: "An extreme close-up of an gray-haired man with a beard in his 60s.".

eration methods. Specifically, for CogVideoX-5B, RainFusion only drops by 0.28% in VBench total score, while DiT-FastAttn and Δ -DiT drop by 5.37% and 5.36% respectively. For OpenSoraPlan-v1.2 and HunyuanVideo, the results is similar that our RainFusion performs best with -0.32% and -0.4% total score loss, respectively.

Qualitative Analysis and Integrability As shown in Fig.5, RainFusion achieves the best visual quality among all methods while DiTFastAttn and Δ -DiT suffer from noise patch or inconsistency subjects. Notably, RainFusion is orthogonal to other acceleration approaches like cached-based method and can be combined to achieve a multiplicative speedup as shown in Fig.1. For example, integrating RainFusion-1.84× with Δ -DiT 1.3× yields a total speedup of 2.4× on HunyuanVideo. As detailed in Tab.2, RainFusion+ variant employs dynamic bandwidth selection (0.5, 0.25, 0.125) across different attention heads. We determine

the optimal bandwidth for each head by selecting the minimum value that maintains 90% recall. Combining RainFusion+ with Δ -DiT results in a minor -0.49% performance drop, demonstrating its practical feasibility.

5. Conclusion

We introduce RainFusion, which utilizes spatial, temporal and textural sparsity in video generation models. Experiments demonstrate that RainFusion can achieve significant speed up in several video generation models with negligible quality loss (~0.2% loss on VBench score). Our method is training-free and calibration-free, making it a plug-and-play tools to speed up video generation models. For future work, we will dive deeper to improve the sparsity ratio while preserving video quality and try to improve the video quality with fine-tuning.

References

- [1] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, Varun Jampani, and Robin Rombach. Stable video diffusion: Scaling latent video diffusion models to large datasets, 2023. 1
- [2] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling latent video diffusion models to large datasets. arXiv preprint arXiv:2311.15127, 2023. 3
- [3] Daniel Bolya and Judy Hoffman. Token merging for fast stable diffusion. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4599–4603, 2023. 3
- [4] Junsong Chen, Chongjian Ge, Enze Xie, Yue Wu, Lewei Yao, Xiaozhe Ren, Zhongdao Wang, Ping Luo, Huchuan Lu, and Zhenguo Li. Pixart-σ: Weak-to-strong training of diffusion transformer for 4k text-to-image generation. In European Conference on Computer Vision, pages 74–91. Springer, 2024. 3
- [5] Pengtao Chen, Mingzhu Shen, Peng Ye, Jianjian Cao, Chongjun Tu, Christos-Savvas Bouganis, Yiren Zhao, and Tao Chen. δ-dit: A training-free acceleration method tailored for diffusion transformers. arXiv preprint arXiv:2406.01125, 2024. 3
- [6] Pengtao Chen, Mingzhu Shen, Peng Ye, Jianjian Cao, Chongjun Tu, Christos-Savvas Bouganis, Yiren Zhao, and Tao Chen. δ-dit: A training-free acceleration method tailored for diffusion transformers, 2024. 2, 5
- [7] Xuanyao Chen, Zhijian Liu, Haotian Tang, Li Yi, Hang Zhao, and Song Han. Sparsevit: Revisiting activation sparsity for efficient high-resolution vision transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2061–2070, 2023. 3
- [8] Prafulla Dhariwal, Jonathan Ho, Ajay Jain, and Pieter Abbeel. Guided diffusion models. In *NeurIPS*, 2022. 3
- [9] Roberto Henschel, Levon Khachatryan, Daniil Hayrapetyan, Hayk Poghosyan, Vahram Tadevosyan, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. Streamingt2v: Consistent, dynamic, and extendable long video generation from text, 2024. 1
- [10] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. Advances in neural information processing systems, 33:6840–6851, 2020. 3
- [11] Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, Yaohui Wang, Xinyuan Chen, Limin Wang, Dahua Lin, Yu Qiao, and Ziwei Liu. Vbench: Comprehensive benchmark suite for video generative models, 2023. 5
- [12] Huiqiang Jiang, Yucheng Li, Chengruidong Zhang, Qianhui Wu, Xufang Luo, Surin Ahn, Zhenhua Han, Amir H. Abdi, Dongsheng Li, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. Minference 1.0: Accelerating pre-filling for long-context llms via dynamic sparse attention, 2024. 3

- [13] Kumara Kahatapitiya, Haozhe Liu, Sen He, Ding Liu, Menglin Jia, Chenyang Zhang, Michael S. Ryoo, and Tian Xie. Adaptive caching for faster video generation with diffusion transformers, 2024. 2, 3
- [14] Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangfeng Xiong, Xin Li, Bo Wu, Jianwei Zhang, Kathrina Wu, Qin Lin, Junkun Yuan, Yanxin Long, Aladdin Wang, Andong Wang, Changlin Li, Duojun Huang, Fang Yang, Hao Tan, Hongmei Wang, Jacob Song, Jiawang Bai, Jianbing Wu, Jinbao Xue, Joey Wang, Kai Wang, Mengyang Liu, Pengyu Li, Shuai Li, Weiyan Wang, Wenqing Yu, Xinchi Deng, Yang Li, Yi Chen, Yutao Cui, Yuanbo Peng, Zhentao Yu, Zhiyu He, Zhiyong Xu, Zixiang Zhou, Zunnan Xu, Yangyu Tao, Qinglin Lu, Songtao Liu, Dax Zhou, Hongfa Wang, Yong Yang, Di Wang, Yuhong Liu, Jie Jiang, and Caesar Zhong. Hunyuanvideo: A systematic framework for large video generative models, 2025. 3, 5
- [15] Lijiang Li, Huixia Li, Xiawu Zheng, Jie Wu, Xuefeng Xiao, Rui Wang, Min Zheng, Xin Pan, Fei Chao, and Rongrong Ji. Autodiffusion: Training-free optimization of time steps and architectures for automated diffusion model acceleration, 2023.
- [16] Feng Liu, Shiwei Zhang, Xiaofeng Wang, Yujie Wei, Haonan Qiu, Yuzhong Zhao, Yingya Zhang, Qixiang Ye, and Fang Wan. Timestep embedding tells: It's time to cache for video diffusion model, 2024. 2, 3
- [17] Haozhe Liu, Wentian Zhang, Jinheng Xie, Francesco Faccio, Mengmeng Xu, Tao Xiang, Mike Zheng Shou, Juan-Manuel Perez-Rua, and Jürgen Schmidhuber. Faster diffusion via temporal attention decomposition. arXiv e-prints, pages arXiv-2404, 2024. 3
- [18] Nanye Ma, Shangyuan Tong, Haolin Jia, Hexiang Hu, Yu-Chuan Su, Mingda Zhang, Xuan Yang, Yandong Li, Tommi Jaakkola, Xuhui Jia, et al. Inference-time scaling for diffusion models beyond scaling denoising steps. arXiv preprint arXiv:2501.09732, 2025. 3
- [19] Xinyin Ma, Gongfan Fang, Michael Bi Mi, and Xinchao Wang. Learning-to-cache: Accelerating diffusion transformer via layer caching. arXiv preprint arXiv:2406.01733, 2024.
- [20] Xinyin Ma, Gongfan Fang, and Xinchao Wang. Deepcache: Accelerating diffusion models for free. In *The IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. 3
- [21] Xin Ma, Yaohui Wang, Gengyun Jia, Xinyuan Chen, Ziwei Liu, Yuan-Fang Li, Cunjian Chen, and Yu Qiao. Latte: Latent diffusion transformer for video generation, 2024. 1
- [22] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International* conference on machine learning, pages 8162–8171. PMLR, 2021. 3
- [23] Open-Sora. Sora prompt, 2024. 5
- [24] OpenAI. Openai. sora, 2024. 1, 3
- [25] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4195–4205, 2023. 3

- [26] Open-Sora Plan. Open-sora plan, 2024. 1, 3, 5
- [27] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023. 3
- [28] Adam Polyak, Amit Zohar, Andrew Brown, Andros Tjandra, Animesh Sinha, Ann Lee, Apoorv Vyas, Bowen Shi, Chih-Yao Ma, Ching-Yao Chuang, David Yan, Dhruv Choudhary, Dingkang Wang, Geet Sethi, Guan Pang, Haoyu Ma, Ishan Misra, Ji Hou, Jialiang Wang, Kiran Jagadeesh, Kunpeng Li, Luxin Zhang, Mannat Singh, Mary Williamson, Matt Le, Matthew Yu, Mitesh Kumar Singh, Peizhao Zhang, Peter Vajda, Quentin Duval, Rohit Girdhar, Roshan Sumbaly, Sai Saketh Rambhatla, Sam Tsai, Samaneh Azadi, Samyak Datta, Sanyuan Chen, Sean Bell, Sharadh Ramaswamy, Shelly Sheynin, Siddharth Bhattacharya, Simran Motwani, Tao Xu, Tianhe Li, Tingbo Hou, Wei-Ning Hsu, Xi Yin, Xiaoliang Dai, Yaniv Taigman, Yaqiao Luo, Yen-Cheng Liu, Yi-Chiao Wu, Yue Zhao, Yuval Kirstain, Zecheng He, Zijian He, Albert Pumarola, Ali Thabet, Artsiom Sanakoyeu, Arun Mallya, Baishan Guo, Boris Araya, Breena Kerr, Carleigh Wood, Ce Liu, Cen Peng, Dimitry Vengertsev, Edgar Schonfeld, Elliot Blanchard, Felix Juefei-Xu, Fraylie Nord, Jeff Liang, John Hoffman, Jonas Kohler, Kaolin Fire, Karthik Sivakumar, Lawrence Chen, Licheng Yu, Luva Gao, Markos Georgopoulos, Rashel Moritz, Sara K. Sampson, Shikai Li, Simone Parmeggiani, Steve Fine, Tara Fowler, Vladan Petrovic, and Yuming Du. Movie gen: A cast of media foundation models, 2024. 1
- [29] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 3
- [30] Amirmojtaba Sabour, Sanja Fidler, and Karsten Kreis. Align your steps: Optimizing sampling schedules in diffusion models, 2024. 2
- [31] Ethan Smith, Nayan Saxena, and Aninda Saha. Todo: Token downsampling for efficient generation of high-resolution images. *arXiv preprint arXiv:2402.13573*, 2024. 3
- [32] Hanlin Tang, Yang Lin, Jing Lin, Qingsen Han, Shikuan Hong, Yiwu Yao, and Gongyi Wang. Razorattention: Efficient kv cache compression through retrieval heads. *arXiv* preprint arXiv:2407.15891, 2024. 3
- [33] Yuchuan Tian, Zhijun Tu, Hanting Chen, Jie Hu, Chao Xu, and Yunhe Wang. U-dits: Downsample tokens in u-shaped diffusion transformers. arXiv preprint arXiv:2405.02730, 2024. 3
- [34] Chau Tran, Duy MH Nguyen, Manh-Duy Nguyen, TrungTin Nguyen, Ngan Le, Pengtao Xie, Daniel Sonntag, James Y Zou, Binh Nguyen, and Mathias Niepert. Accelerating transformers with spectrum-preserving token merging. Advances in Neural Information Processing Systems, 37:30772–30810, 2025. 3
- [35] Hongjie Wang, Difan Liu, Yan Kang, Yijun Li, Zhe Lin, Niraj K Jha, and Yuchen Liu. Attention-driven training-free efficiency enhancement of diffusion models. In *Proceedings*

- of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 16080–16089, 2024. 3
- [36] Felix Wimbauer, Bichen Wu, Edgar Schoenfeld, Xiaoliang Dai, Ji Hou, Zijian He, Artsiom Sanakoyeu, Peizhao Zhang, Sam Tsai, Jonas Kohler, et al. Cache me if you can: Accelerating diffusion models through block caching. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 6211–6220, 2024. 3
- [37] Haoyu Wu, Jingyi Xu, Hieu Le, and Dimitris Samaras. Importance-based token merging for diffusion models. arXiv preprint arXiv:2411.16720, 2024. 3
- [38] Haocheng Xi, Shuo Yang, Yilong Zhao, Chenfeng Xu, Muyang Li, Xiuyu Li, Yujun Lin, Han Cai, Jintao Zhang, Dacheng Li, et al. Sparse videogen: Accelerating video diffusion transformers with spatial-temporal sparsity. arXiv preprint arXiv:2502.01776, 2025. 2, 4, 6
- [39] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, Da Yin, Xiaotao Gu, Yuxuan Zhang, Weihan Wang, Yean Cheng, Ting Liu, Bin Xu, Yuxiao Dong, and Jie Tang. Cogvideox: Text-to-video diffusion models with an expert transformer, 2024. 3, 5
- [40] Zhihang Yuan, Hanling Zhang, Pu Lu, Xuefei Ning, Linfeng Zhang, Tianchen Zhao, Shengen Yan, Guohao Dai, and Yu Wang. Ditfastattn: Attention compression for diffusion transformer models, 2024. 2, 3, 5
- [41] Stephen Zhang and Vardan Papyan. Oats: Outlier-aware pruning through sparse and low rank decomposition. *arXiv* preprint arXiv:2409.13652, 2024. 3
- [42] Zangwei Zheng, Xiangyu Peng, Tianji Yang, Chenhui Shen, Shenggui Li, Hongxin Liu, Yukun Zhou, Tianyi Li, and Yang You. Open-sora: Democratizing efficient video production for all. arXiv preprint arXiv:2412.20404, 2024. 3