

---

# Fast and Memory-Efficient Video Diffusion Using Streamlined Inference

---

Zheng Zhan<sup>1\*</sup> Yushu Wu<sup>1\*</sup> Yifan Gong<sup>1</sup> Zichong Meng<sup>1</sup> Zhenglun Kong<sup>1,2</sup>  
Changdi Yang<sup>1</sup> Geng Yuan<sup>3</sup> Pu Zhao<sup>1†</sup> Wei Niu<sup>3</sup> Yanzhi Wang<sup>1</sup>  
<sup>1</sup>Northeastern University <sup>2</sup>Harvard University <sup>3</sup>University of Georgia

512×768 resolution

Animatediff  
Peak Mem: 29.7G  
Latency: 16.9s

Ours  
Peak Mem:  
9.08G (3.3× less)  
Latency:  
10.7s (1.6× less)



*a full-sized man rides a comically small motorcycle through a residential neighborhood*

Figure 1: Our Streamlined Inference is a training-free inference framework for video diffusion models that can reduce the computation and peak memory cost without sacrificing the quality.

## Abstract

The rapid progress in artificial intelligence-generated content (AIGC), especially with diffusion models, has significantly advanced development of high-quality video generation. However, current video diffusion models exhibit demanding computational requirements and high peak memory usage, especially for generating longer and higher-resolution videos. These limitations greatly hinder the practical application of video diffusion models on standard hardware platforms. To tackle this issue, we present a novel, training-free framework named Streamlined Inference, which leverages the temporal and spatial properties of video diffusion models. Our approach integrates three core components: Feature Slicer, Operator Grouping, and Step Rehash. Specifically, Feature Slicer effectively partitions input features into sub-features and Operator Grouping processes each sub-feature with a group of consecutive operators, resulting in significant memory reduction without sacrificing the quality or speed. Step Rehash further exploits the similarity between adjacent steps in diffusion, and accelerates inference through skipping unnecessary steps. Extensive experiments demonstrate that our approach significantly reduces peak memory and computational overhead, making it feasible to generate high-quality videos on a single consumer GPU (e.g., reducing peak memory of AnimateDiff from 42GB to 11GB, featuring faster inference on 2080Ti)<sup>1</sup>.

## 1 Introduction

Recent years have witnessed continual progress and advancements in artificial intelligence-generated content (AIGC). Among them, diffusion models allow artists and amateurs to create visual content

---

\*Equal contributions.

†Corresponding author

<sup>1</sup>Code available at: <https://github.com/wuyushuwys/FMEDiffusion>

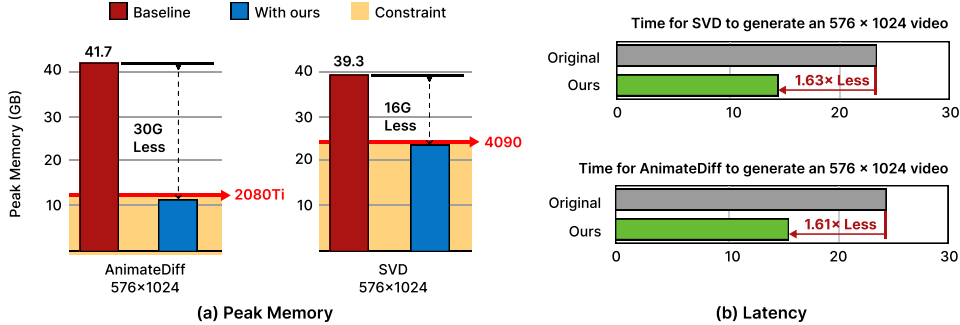


Figure 2: Comparison on Animatediff and SVD inference using our Streamlined Inference. Memory requirement is crucial as “**Out of Memory**” errors prevent the GPU from performing inference.

with text prompts, advancing the development of image and video generation in both academia and industry. For video diffusion models, the latest works such as SVD-XT [3], Gen2 [29], Pika [20], and notably the more advanced Sora [28], demonstrate impressive capabilities in producing visually striking and artistically effective videos. Despite their great performance, video diffusion models also exhibit high computational requirements and substantial peak memory, particularly when generating longer videos with higher resolutions. For instance, SVD-XT generates 25 frames simultaneously with a resolution of  $576 \times 1024$ , while Sora expands these capabilities by supporting the generation of longer videos (over a minute) at a higher resolution of  $1080 \times 1920$ . Given the trends of generating longer videos with higher quality, the escalating memory and computation demands have impeded practical applications of these large-scale video diffusion models on various platforms.

Existing model compression methods to reduce peak memory and latency, such as weight pruning [38, 12, 35, 6, 48, 23, 42, 44, 47], quantization [37, 22, 32], and distillation [18, 10, 43, 39], typically require substantial retraining or fine-tuning of the compressed model to recover performance. This process is costly, time-consuming, and may raise data privacy concerns. Applying these methods in zero-shot avoids the expensive retraining, but leads to severe performance degradation. Furthermore, the variety and complexity of video diffusion architectures further complicate the model optimization. Therefore, it is challenging yet crucial to develop an effective and efficient video diffusion framework with reduced computations, smaller peak memory and less data (no re-training) requirements for its wide applications.

To address the above challenges, we first identify the sources of the **high computation and memory cost**, which scale up with the iterative denoising process and the simultaneous processing of multiple frames. We further observe that the feature maps of certain layers may exhibit high similarity between multiple consecutive denoising steps due to the temporal property of videos, enabling further optimizations for acceleration. Based on that, we propose a training-free framework named Streamlined Inference, by leveraging the temporal and spatial characteristics of video diffusion models to effectively reduce peak memory and computational demands. Our framework contains three core components: Feature Slicer, Operator Grouping, and Step Rehash, which work together closely and comprehensively with different focuses on peak memory reduction or inference acceleration.

Our Feature Slicer performs lossless feature slicing in both temporal and spatial layers, raising the possibility of peak memory reduction through processing smaller features. However, the feature slicer alone is not able to decrease peak memory as we still need to store all intermediate results of one layer for all sliced features to form a complete intermediate feature map for the next layer. To reduce peak memory practically, we further propose Operator Grouping to group homogeneous and consecutive operators in the computational graph. Within each operator group, the intermediate result of one sliced feature can be directly sent to the next operator/layer without waiting for aggregation with all other intermediate results, achieving the full potential of Feature Slicer to reduce the peak memory. Furthermore, a pipeline technique is proposed to accelerate the computations of the same operator group for multiple sliced feature inputs, with improved parallelism.

Moreover, observing the high similarity of certain features between multiple consecutive denoising steps, we propose Step Rehash to reuse the generated features for a few following steps due to their high similarity, skipping the exact expansive and repetitive generation of similar features and thereby accelerating the video diffusion significantly. With this framework, we can generate high-quality

videos in a fast and memory-efficient manner on a single consumer GPU, as shown in Fig. 2. For example, the peak memory of AnimateDiff [11] can be reduced significantly from 41.7GB to 11GB, featuring inference on a typical consumer GPU 2080Ti. We summarize our contributions as follows:

- We propose a novel training-free framework that can significantly reduce the peak memory and computation cost for the inference of video diffusion models by leveraging the spatial and temporal characteristics of video diffusion models.
- Our approach can be seamlessly integrated into existing video diffusion models. Our extensive experiments on SVD, SVD-XT, and AnimateDiff demonstrate our effectiveness to reduce peak memory and accelerate inference without sacrificing quality.
- Our approach offers a new research perspective for fast and memory-efficient video diffusion, enabling the generation of higher quality and longer videos on consumer-grade GPUs.

## 2 Related Work

**Video Diffusion Models.** For video generation, various approaches have been proposed, with VDM [17] as a leading example. VDM transforms the conventional U-Net [30] architecture of image diffusion models into a 3D U-Net structure, employing joint training on both images and videos. MagicVideo [49] is the first work that introduces Latent Diffusion Model (LDM) for text-to-video (T2V) generation in latent space. LVDM [13] introduces a mask sampling technique that enhances its longer video generation capability. ModelScope [36] incorporates spatial-temporal convolution and attention into LDM. Video LDM [4] trains a T2V network composed of three training stages, enabling higher quality and longer video generation. Show-1 [45] first introduces the fusion of pixel-based and latent-based diffusion models for T2V generation. Recently, Stable Video Diffusion (SVD) [3] identifies three key stages for training video LDMs: text-to-image (T2I) pretraining, video pretraining, and high-quality video finetuning.

**Architectural Efficiency of Video Diffusion Models.** There are various research efforts exploring either architectural efficiency or model compression techniques for image/video generation. For example, ED-T2V [24] freezes parameters to reduce training costs and proposes an attention mechanism to ensure temporal coherence. SimDA [40] devises a parameter-efficient training approach by maintaining the parameter of the T2I model and uses two adapters to train it. For model compression, Diff-pruning [6] employs structural pruning techniques to reduce inference time at each sampling step. Additionally, the work [22] implements quantization on diffusion models using low-precision data types. However, these methods take substantial efforts to retrain or finetune the diffusion model to recover performance, which is costly, time-consuming, and may raise data privacy concerns. Furthermore, applying post-training compression techniques in one-shot [8, 7, 34] may save the retraining/fine-tuning efforts, but suffers from significant performance degradation.

**Sample-Efficient Video Diffusion Models.** To address the iterative denoising process in diffusion models and improve the sampling efficiency, two approaches are proposed. The first approach [2, 19, 5, 25] focuses on creating rapid solvers to resolve the differential equation associated with the denoising process more effectively. The works [31, 27, 21] utilize knowledge distillation methods to compress and simplify the sampling trajectory efficiently, thereby enhancing overall performance. Imagen video [15] is one of the first methods to apply progressive distillation on video diffusion models, incorporating guidance and stochastic samplers. Recent work Deepcache [26] proposes a novel training-free paradigm that accelerates diffusion models by reusing the high-level features.

## 3 Motivation

**Peak memory and computation analysis.** Existing open-source video diffusion models [11, 3, 36, 46] typically adopt a pretrained T2I 2D-UNet as backbone. Their temporal layers are seamlessly integrated into the backbone 2D-UNet, positioned after every spatial layer. Here, we use SVD as an example to demonstrate how peak memory and computational overhead scale with the number of frames. The SVD model is trained with two distinct configurations: regular SVD is designed to generate 14 frames, while SVD-XT is tailored to produce 25 frames. To generate 14 or 25 video frames concurrently with SVD, its latent features require massive GPU memory and computation consumption, estimated to be approximately  $14\times$  or  $25\times$  higher than its base T2I model. This



Figure 3: The quality results of our method and naïve slicing. Note that naïve slicing will incur unpleasant artifacts due to lack of temporal correction by fewer frames.

estimate does not even account for the additional memory required by SVD’s extra-temporal layers. More specifically, the SVD consumes 39.49G of peak memory for  $576 \times 1024$  resolution output, whereas its image generation counterpart only requires 6.33G of memory at the same resolution. Furthermore, incorporating the classifier-free guidance [16] substantially enhances the generation quality but doubles the peak memory required during inference.

Consequently, video diffusion is computationally demanding, but the challenge of memory consumption is more critical and demands immediate attention. Most consumer-grade GPUs do not have enough memory for video diffusion models and, therefore, suffer from the “**Out of Memory**” error, which prevents the GPU from generating high-quality videos. There is no workaround without switching GPUs. Most users have to endure generating short and low-resolution videos.

**Naïve Slicing.** A Naïve approach to reduce peak memory is to execute the video diffusion inference clip-by-clip. However, this strategy is hindered by the temporal layers, which are essential for maintaining temporal correlation in video diffusion models. Forcibly implementing this approach can generate random artifacts and cause motion vanishing in the output video, as detailed in Fig. 3. Therefore, designing a memory-efficient inference framework is a challenging and non-trivial task.

## 4 Streamlined Inference Framework

To address the above massive peak memory and computation costs, in this section, we propose a training-free framework named Streamlined Inference, which is composed of three core components: Feature Slicer, Operator Grouping, and Step Rehash. First, we discuss Feature Slicer, designed to partition input features of spatial and temporal layers, and enable the potential of massive peak memory reduction. Next, we introduce our Operator Grouping technique to aggregate homogeneous and consecutive operators into the same group, achieving the full potential of Feature Slicer to reduce peak memory through reusing the memory of intermediate result from previous sliced feature. Finally, we discuss our Step Rehash method to reuse the same feature for a few consecutive steps due to their high similarity. It accelerates the inference without increasing peak memory overhead as it skips certain denoising steps with less computations.

### 4.1 Feature Slicer

Video diffusion models contain spatial and temporal layers which extract the corresponding information from their specific domains. On this basis, we propose a feature slicer that consists of two components: Spatial-layer slicer and Temporal-layer slicer, to divide the feature map into multiple batches/sub-features, ensuring accurate computation without introducing additional operations. The slicer is further utilized for Operator Grouping to reduce peak memory cost.

**Spatial layers slicer.** Based on our profiling (more details can be found in Appendix A) for memory allocation of various video diffusion models, we find that performing slices at spatial layers can greatly reduce the memory footprint. The 5-D feature in the spatial layer  $X \in \mathbb{R}^{B \times T \times C \times H \times W}$  can be reshaped to a 4-D feature  $X \in \mathbb{R}^{(B \times T) \times C \times H \times W}$ , where  $B, T, C, H, W$  are the batch size, number of frames, channels, height, and width, respectively. Thus, we slice it into  $k$  sub-features,  $\{X_i \in \mathbb{R}^{n_i \times C \times H \times W}\}_{i=1}^k$ , where  $n_i = \lceil B \times T/k \rceil$  with  $\lceil \cdot \rceil$  denoting the least integer greater than or equal to the input. If  $\lceil B \times T/k \rceil \neq B \times T/k$ , the dimension of the last sub-feature  $n_k$  is different

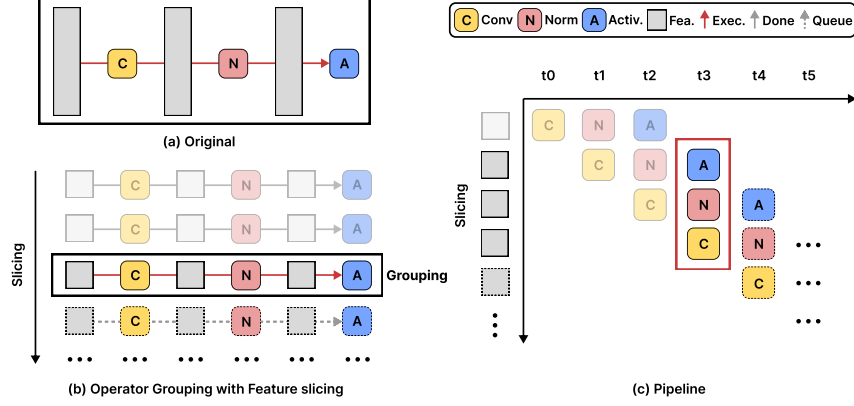


Figure 4: Overview of Operator Grouping with Pipeline in our framework.

from others. The spatial layer slicer is applicable for most operations in spatial layers such as Conv2D, GroupNorm, LayerNorm, Attention, and Linear.

**Temporal layers slicer.** The input of the temporal layer is a 5-D feature map with dimensions  $\{batch, channels, frames, height, width\}$ . 3-D operations such as Conv3D are employed to extract temporal information from the 5-D feature. Differing from spatial layers, slicing along the temporal dimension may disrupt the extraction and processing of temporal information. Therefore, we keep the temporal dimension untouched while slicing over other dimensions. Specifically, the 5-D feature  $X \in \mathbb{R}^{B \times T \times C \times H \times W}$  can be sliced to  $k_h \times k_w$  sub-features  $\{X_{ij} \in \mathbb{R}^{B \times T \times C \times h_i \times w_j}\}_{i=1, j=1}^{i=k_h, j=k_w}$ , where  $h_i = \lceil H/k_h \rceil$  and  $w_j = \lceil W/k_w \rceil$ . After detailed profiling different configurations for temporal layer slicer, we discover that the configuration with  $k_h = \max(H, 16)$  and  $k_w = \max(W, 16)$  can result in promising peak memory reduction.

## 4.2 Operator Grouping

Although Feature Slicer converts the original feature map into multiple smaller sub-features with reduced memory footprint, the peak memory can not be reduced since the intermediate results of multiple sliced features require re-consolidation to send to the next layer/operator as inputs. It still needs to store all intermediate outputs of sliced features to form the united/unsliced intermediate feature map without actual peak memory reduction. Therefore, to address this problem, we propose Operator Grouping to group the operators accordingly in the computational graph, achieving the full potential of Feature Slicer with effective peak memory reduction due to less memory reserved for intermediate results. Furthermore, a pipeline technique is proposed to optimize the inference of operator groups with improved parallelism and practical acceleration.

### 4.2.1 Grouping Operators for Peak Memory Reduction

Operator Grouping directly re-uses existing operators by aggregating consecutive homogeneous operators into the same group. Homogeneous operators indicate these operators extract features from coherent domains and dimensions. In video diffusion models, different operators can be grouped into  $\text{GroupOP}^t$  (temporal operator groups) and  $\text{GroupOP}^s$  (spatial operators groups) to ensure the well-preserved semantics of sliced sub-features within each group. For example, in the SVD Model [3], consecutive GroupNorm, Conv2D, SiLU, and Up/DownSample operators in the *Spatial ResBlock* can be aggregated to one group, as these operators all extract features from spatial domain and are deemed homogeneous. As shown in Fig. 4, when computing the output feature  $X^o$  for an operator group (GroupOP), the input feature  $X$  is sliced into multiple sub-features  $X_1, X_2, \dots, X_k$  with Feature Slicer. Each sub-feature  $X_i$  goes through the operator group and their outputs are concatenated after all outputs are available, as shown in Eq. (1),

$$X^o = \text{Concat}(\text{GroupOP}(X_1), \text{GroupOP}(X_2), \dots, \text{GroupOP}(X_k)) \quad (1)$$

where  $k$  is the number of slices, and Concat is the concatenation operation.

**Reducing peak memory cost.** As shown in Fig. 4, the peak memory with the operator group is determined by the memory footprint of the input feature, the output feature, and the intermediate

results. Without operator grouping, all intermediate results of all operators for sliced sub-features will allocate their own memory, hence failing to reduce memory consumption. Compared with the case above, grouped operators only need to allocate memory for intermediate results of a single sliced sub-feature and the final outputs, without the necessity to allocate full intermediate features corresponding to the original unsliced input feature, as shown in Fig. 4 (a) and (b). Operator Grouping can effectively reduce peak memory cost, enabling successful inference of video diffusion models on one single consumer or commercial GPU with low or moderate available memory, as shown in Tab. 1.

**Mitigating I/O intensity.** As the original feature map is sliced into multiple sub-features to reduce peak memory cost, the computation may slow down due to multiple iterations corresponding to multiple inputs. However, we surprisingly observe that even with the naive basic for-loop implementation for each sub-feature as shown in Fig. 4 (b), the overall runtime with Operator Grouping is around 10% slower than that of the original unsliced version. The marginal increase in runtime can be attributed to the memory bound of the GPU for video diffusion inference. Specifically, current video diffusion model inference suffers from the memory bound, where the I/O overhead of intermediate results is more notorious than their computation workload. The slicer provides a solution to mitigate the I/O burden, thus balancing the computation and memory read/write to fully utilize the GPU capacity.

#### 4.2.2 Pipelining with Improved Parallelism and Practical Acceleration.

With the proposed Feature Slicer and Operator Grouping, the peak memory will decrease significantly with a marginal increase for the computation runtime (based on the basic for-loop implementation). With a deeper investigation for the computation patterns, we find that the *for-loop implementation* cannot maximize GPU parallelism, and further employ the pipelining technique to optimize the *for-loop implementation* for faster inference without additional memory cost.

With Operator Grouping, there are multiple operators in one group to process one sliced sub-feature sequentially. With the naive *for-loop implementation*, before feeding each sliced sub-feature into another operator group, it needs to wait until the last sub-feature finishes its computation within the group. The parallelism can be further improved with the proposed pipeline method. Specifically, in an operator group, after a sliced feature map is computed by the out-of-place computation operator (e.g., Conv, GroupNorm, Attention, etc.) and sent to the next operator, its previous allocated memory is no longer required, but it is still reserved during inference, leading to resource waste. We can pipeline all operators in the same group to mitigate this issue. As shown in Fig. 4(c), once the Conv operator completes processing a sliced feature  $X_i$  as described in Eq. (1) and its outputs are sent to the next operator Norm, the subsequent sliced feature  $X_{i+1}$  is immediately piped into the same Conv operator, reusing the reserved memory of  $X_i$ . In this way, multiple operators are executed simultaneously with improved parallelism. No additional memory is required, as we only make use of previously reserved memory.

**Acceleration performance.** With the naive *for-loop implementation*, only one operator in an operator group is executed at a time. However, our pipeline method can simultaneously execute multiple operators (such as Conv, Norm, and Activation as depicted in Fig. 4 (c)) without incurring additional memory. Consequently, the inference speed can be further improved. Accordingly, integrating the pipeline within Operator Grouping can mitigate 10% speed degradation caused by feature slicing.

### 4.3 Step Rehash

In this section, we further introduce our step rehash method to optimize the iterative denoising steps for effective acceleration in video diffusion generation. Capitalized on the high similarity between adjacent steps, our approach accelerates the video generation, while ensuring both high quality and temporal consistency across video frames, without extra memory cost. Next we first discuss our observations for the high feature similarity and then explain details of our step rehash.

#### 4.3.1 Similarity of Temporal Features between Steps

**Similarity visualization.** The denoising process of U-Net in diffusion models requires multiple steps and the features of different steps may share certain similarities with minor differences [26]. To explore this, we analyze the feature maps averaged over multiple images at different parts of the model and plot the similarity between features of different steps, with an example shown in Fig. 5 (and more results and details demonstrated in Appendix C). We find two key insights below:

- The similarity between adjacent steps significantly depends on certain blocks and layers, and it may change sharply after specific operations in video diffusion. The features do not always show high similarity. For example, neither deeper layers within the same block nor those in middle blocks consistently show higher similarity between adjacent steps.
- The features between adjacent steps following the temporal layers and spatial layers in video diffusion usually exhibit remarkably higher similarity compared to outputs of other layers.

**High similarity after temporal layers.** Existing video diffusion models typically employ pretrained image diffusion models as their backbone. While these image models are trained to produce a variety of images, the addition of temporal layers is designed to improve the temporal continuity of latent features. This enhancement significantly strengthens their correlation, thereby increasing similarity among the features.

**Motivation and challenges for step rehash.**

Due to the high similarity between features of different steps, we propose the step rehash method to skip the computation of certain features by reusing previous generated features. However, we need to address the challenges of when and where to skip. Specifically, based on the above insights, simply reusing features from deeper layers does not guarantee better results since deeper layers may not show high similarity. We need to carefully choose what layers or blocks can be skipped (**where to skip**) to make use of high similarities without significantly downgrading the generation performance. As shown in fig. 5, it exhibits high similarity between adjacent time steps, but the similarity pattern differs between video diffusion models. Thus, we need to determine which steps can use skip strategy (**when to skip**), and the remaining steps that require full computation are *full computation steps*.

**4.3.2 Step Rehash**

Here we specify the details of our step rehash. The video diffusion models typically use a U-Net architecture with 4 down-sampling and 4 up-sampling blocks, and their output features can be represented by  $D_{0\sim 3}^s$  and  $U_{0\sim 3}^s$ , respectively, with  $s$  denoting the current step number as shown in Fig. 6. Typically,  $U_b^s$  is obtained by feeding  $D_b^s$  and  $U_{b-1}^s$  into the  $b^{th}$  ( $b > 0$ ) up-sampling block, and  $U_3^s$  is the final output of the  $s^{th}$  step. Based on similarity analysis, in the next step  $s + 1$ , we can directly reuse the output features of the temporal layer from the previous step  $s$  without actual exact computations. Our insights into the similarity indicate that deeper and middle blocks do not consistently demonstrate high similarity. Reusing their features results in significant degradation of generation quality. Therefore, we rehash features of the temporal layer in the final up-sampling block. Specifically, to obtain  $U_3^{s+1}$  for step  $s + 1$ , we feed the output features of the temporal layer from  $U_3^s$  (current *full computation step*) into the final up-sampling block. Note that we only compute part of  $U_3^{s+1}$  and do not need to compute  $D_3^{s+1}$  for concatenation, since our reused temporal layer is deeper than the concat operator for features from  $D_3^{s+1}$ , as shown in Fig. 6. We further propose a step search algorithm to solve the **when to skip** problem, algorithm details can be found in Appendix B.

**5 Experimental results**

**5.1 Models, Datasets and Evaluation Metrics**

We conduct the experiments on representative video diffusion models, including SVD [3], SVD-XT [3], and AnimateDiff [11]. For evaluation, we use the following evaluation protocols: The first frame of the video clips are extracted as the image condition for image-to-video generation and their captions are considered as the prompts. All experiments are conducted on a NVIDIA A100 GPU.

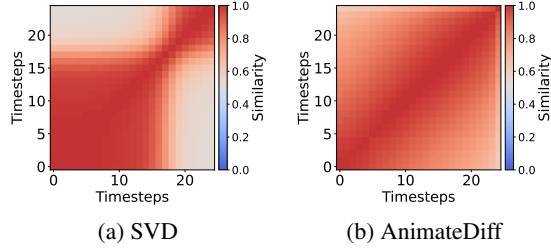


Figure 5: The high similarity of output features after temporal layers in  $U_3$  between each timestep.

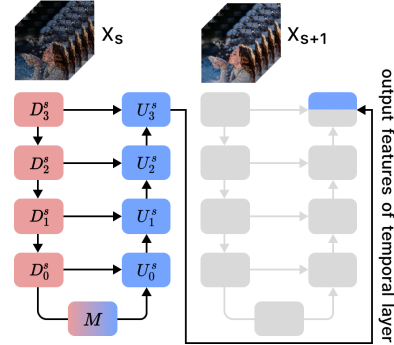


Figure 6: Illustration of Step Rehash. Computation in grey areas are skipped.

Table 1: Comparison of our Streamlined Inference with baseline methods in video visual quality (on UCF101), PM (Peak Memory), and latency (measured with 50 runs with the average value).

Model	Method	FVD↓	CLIP-Score↑	512 × 512		576 × 1024	
				PM	Latency	PM	Latency
SVD #F=14	Original	307.7	29.25	20.91G	10.23s	39.49G	23.29s
	Naïve Slicing	1127.5	26.32	8.12G	31.85s	10.72G	65.56s
	<b>Ours</b>	340.6	28.98	13.67G	7.36s	23.42G	14.24s
SVD-XT #F=25	Original	387.9	28.18	31.97G	17.05s	61.17G	40.77s
	Naïve Slicing	2180.0	24.42	8.12G	59.86s	10.72G	121.82s
	<b>Ours</b>	424.7	27.94	19.37G	12.10s	36.32G	25.47s
AnimateDiff #F=16	Original	758.7	28.89	21.83G	9.65s	41.71G	24.38s
	Naïve Slicing	2403.9	26.63	7.22G	19.98s	9.92G	38.69s
	<b>Ours</b>	784.5	28.71	7.51G	7.08s	11.07G	15.15s

**Zero-shot UCF-101 [33]:** We sample clips from each categories of UCF-101 dataset, and gather a subset with 1,000 video clips for evaluation. Their action categories are considered as their captions. For SVD and SVD-XT, our samples are generated at a resolution of  $576 \times 1024$  (14 frames for SVD and 25 frames for SVD-XT) and then resize to  $240 \times 320$ . For AnimateDiff, we generate samples with resolution  $512 \times 512$  (16 frames).

**Zero-shot MSR-VTT [41]:** We generated a video sample for each of the 9,940 development prompts. The samples are at resolution  $320 \times 576$  then resized to  $240 \times 426$  for all models with different number of generated frames.

**Metrics:** We compute the FVD as outlined in [9] and CLIP-Score [14] using TorchMetrics [1] to measure the performance of generated samples.

**Baseline:** We use pretrained weight for SVD (I2V) and AnimateDiff (T2V). We compare the proposed Streamlined Inference (use 13 *full computation steps*) with the original inference (use 25 *full computation steps*) and naïve slicing inference as mentioned in Sec.3. More specifically, for image-conditioned SVD model, we set each naïve slice with a frame size of 2 and use the last frames of each generated slice as the image condition for the next slice. For AnimateDiff, we evenly generate 4 slices with a frame size of 4, then combine them into a full video clip.

## 5.2 Quantitative Evaluation

The results from Table 1 demonstrate the effectiveness of our proposed method in managing memory, computational resources, and performance. Our method significantly reduced peak memory and latency while maintaining competitive FVD and CLIP-Score values across all three models and resolutions compared to the original method. For SVD, our method achieved a notable reduction in peak memory and latency while maintaining competitive FVD and CLIP-Score, unlike Naïve Slicing, which increased FVD and latency. For SVD-XT, our method improved over Naïve Slicing and balanced resource usage and performance. For AnimateDiff, our method significantly outperformed Naïve Slicing in FVD and latency, achieving nearly the same performance as the original method but with smaller latency and around a 70% reduction in peak memory.

## 5.3 Ablation Study

Our ablation study in Table 2 demonstrates that our Step Rehash method consistently outperforms DeepCache with the same number of *full computation steps*. Step Rehash skips more computations than DeepCache. For the SVD model, our method maintains competitive CLIP-Scores while slightly increasing FVD compared to the original method (FVD of 307.7 and CLIP-Score of 29.25 on UCF101). DeepCache performs poorly, increasing FVD and reducing video quality. For the AnimateDiff model, our method maintains stable FVD (603.9 vs. 607.13) and CLIP-Score (29.29 vs. 29.40) on MSR-VTT compared to the original method. DeepCache shows the worst performance on UCF101, with higher FVD and lower CLIP-Scores. Visual comparisons of our method with DeepCache are provided in Appendix D.



Table 2: Ablation study of our proposed method compared with DeepCache in video visual quality. Both our Step Rehash and DeepCache involve 13 *full computation steps*.

Model	Method	UCF101		MSR-VTT	
		FVD↓	CLIP-Score↑	FVD↓	CLIP-Score↑
SVD	Original	307.7	29.25	373.6	26.06
	DeepCache	394.0	28.57	463.6	25.30
	<b>Step Rehash</b>	340.6	28.98	402.1	25.86
AnimateDiff	Original	758.7	28.89	607.1	29.40
	DeepCache	840.2	28.15	615.8	29.06
	<b>Step Rehash</b>	784.5	28.71	603.9	29.29

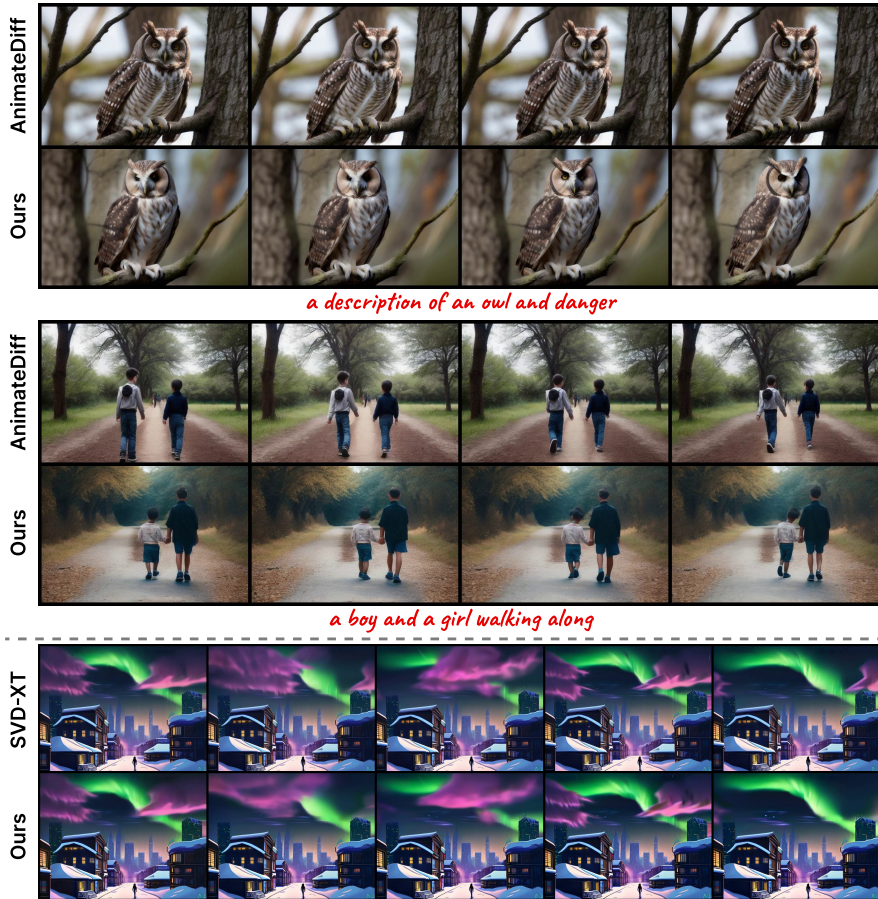


Figure 7: Quality evaluation of using our method on baseline models. The results show that our method can be generally applied to various video diffusion models and achieve competitive results.

#### 5.4 Quality results

In Fig. 7 and Appendix E, we present qualitative results comparing our method to the original model without using Streamlined Inference. We can see that our method produces vivid and high-quality samples aligned with the text descriptions. More importantly, these results demonstrate that our method can significantly reduce peak memory and computation without sacrificing quality.

## 6 Conclusion and Limitation

In this paper, we propose a novel training-free framework that significantly reduces peak memory and computation costs for video diffusion model inference by leveraging its spatial and temporal characteristics. Our approach can be seamlessly integrated into existing models. Extensive experiments on

SVD, SVD-XT, and AnimateDiff demonstrate our method’s effectiveness in reducing peak memory and accelerating inference without sacrificing quality. Our approach offers a new perspective for fast, memory-efficient video diffusion, enabling the generation of higher quality and longer videos on consumer-grade GPUs. Though our method is general, the efficiency is limited by baseline model architecture design.

## Acknowledgments and Disclosure of Funding

This work is supported by National Science Foundation CNS-2312158. We would like to express our sincere gratitude to the reviewers for their invaluable feedback and constructive comments to improve the paper.

## References

- [1] Torchmetrics, [https://lightning.ai/docs/torchmetrics/stable/multimodal/clip\\_score.html](https://lightning.ai/docs/torchmetrics/stable/multimodal/clip_score.html).
- [2] Fan Bao, Chongxuan Li, Jun Zhu, and Bo Zhang. Analytic-dpm: an analytic estimate of the optimal reverse variance in diffusion probabilistic models. *arXiv preprint arXiv:2201.06503*, 2022.
- [3] 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.
- [4] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22563–22575, 2023.
- [5] Tim Dockhorn, Arash Vahdat, and Karsten Kreis. Genie: Higher-order denoising diffusion solvers. *Advances in Neural Information Processing Systems*, 35:30150–30166, 2022.
- [6] Gongfan Fang, Xinyin Ma, and Xinchao Wang. Structural pruning for diffusion models. In *Advances in Neural Information Processing Systems*, 2023.
- [7] Elias Frantar and Dan Alistarh. Sparsegpt: Massive language models can be accurately pruned in one-shot. In *International Conference on Machine Learning*, pages 10323–10337. PMLR, 2023.
- [8] Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. Gptq: Accurate post-training compression for generative pretrained transformers. *arXiv preprint arXiv:2210.17323*, 1, 2022.
- [9] Songwei Ge, Aniruddha Mahapatra, Gaurav Parmar, Jun-Yan Zhu, and Jia-Bin Huang. On the content bias in fr chet video distance. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024.
- [10] Yifan Gong, Zheng Zhan, Qing Jin, Yanyu Li, Yerlan Idelbayev, Xian Liu, Andrey Zharkov, Kfir Aberman, Sergey Tulyakov, Yanzhi Wang, et al. E<sup>2</sup>gan: Efficient training of efficient gans for image-to-image translation. *arXiv preprint arXiv:2401.06127*, 2024.
- [11] Yuwei Guo, Ceyuan Yang, Anyi Rao, Yaohui Wang, Yu Qiao, Dahua Lin, and Bo Dai. AnimateDiff: Animate your personalized text-to-image diffusion models without specific tuning. *arXiv preprint arXiv:2307.04725*, 2023.
- [12] Yang He, Ping Liu, Ziwei Wang, Zhilan Hu, and Yi Yang. Filter pruning via geometric median for deep convolutional neural networks acceleration. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [13] Yingqing He, Tianyu Yang, Yong Zhang, Ying Shan, and Qifeng Chen. Latent video diffusion models for high-fidelity long video generation. *arXiv preprint arXiv:2211.13221*, 2022.

- [14] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7514–7528, 2021.
- [15] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022.
- [16] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*, 2021.
- [17] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *Advances in Neural Information Processing Systems*, 35:8633–8646, 2022.
- [18] Minguk Kang, Richard Zhang, Connelly Barnes, Sylvain Paris, Suha Kwak, Jaesik Park, Eli Shechtman, Jun-Yan Zhu, and Taesung Park. Distilling diffusion models into conditional gans. *arXiv preprint arXiv:2405.05967*, 2024.
- [19] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. *Advances in Neural Information Processing Systems*, 35:26565–26577, 2022.
- [20] Pika Labs. Pika labs ai text to video generator, 2023.
- [21] Muyang Li, Ji Lin, Chenlin Meng, Stefano Ermon, Song Han, and Jun-Yan Zhu. Efficient spatially sparse inference for conditional gans and diffusion models. *arXiv preprint arXiv:2211.02048*, 2022.
- [22] Xiuyu Li, Yijiang Liu, Long Lian, Huanrui Yang, Zhen Dong, Daniel Kang, Shanghang Zhang, and Kurt Keutzer. Q-diffusion: Quantizing diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 17535–17545, 2023.
- [23] Yanyu Li, Pu Zhao, Geng Yuan, Xue Lin, Yanzhi Wang, and Xin Chen. Pruning-as-search: Efficient neural architecture search via channel pruning and structural reparameterization. *arXiv preprint arXiv:2206.01198*, 2022.
- [24] Jiawei Liu, Weining Wang, Wei Liu, Qian He, and Jing Liu. Ed-t2v: An efficient training framework for diffusion-based text-to-video generation. In *2023 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8, 2023.
- [25] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps. *Advances in Neural Information Processing Systems*, 35:5775–5787, 2022.
- [26] 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.
- [27] Chenlin Meng, Ruiqi Gao, Diederik P Kingma, Stefano Ermon, Jonathan Ho, and Tim Salimans. On distillation of guided diffusion models. *arXiv preprint arXiv:2210.03142*, 2022.
- [28] OpenAI. Creating video from text, 2024.
- [29] Runway Research. Gen-2: The next step forward for generative ai, 2023.
- [30] 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.
- [31] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. *arXiv preprint arXiv:2202.00512*, 2022.
- [32] Xuan Shen, Zhenglun Kong, Changdi Yang, Zhaoyang Han, Lei Lu, Peiyan Dong, et al. Edgeqat: Entropy and distribution guided quantization-aware training for the acceleration of lightweight lms on the edge. *arXiv preprint arXiv:2402.10787*, 2024.

- [33] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012.
- [34] Mingjie Sun, Zhuang Liu, Anna Bair, and J Zico Kolter. A simple and effective pruning approach for large language models. *arXiv preprint arXiv:2306.11695*, 2023.
- [35] Haotao Wang, Shupeng Gui, Haichuan Yang, Ji Liu, and Zhangyang Wang. Gan slimming: All-in-one gan compression by a unified optimization framework. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IV 16*, pages 54–73. Springer, 2020.
- [36] Jiuniu Wang, Hangjie Yuan, Dayou Chen, Yingya Zhang, Xiang Wang, and Shiwei Zhang. Modelscope text-to-video technical report. *arXiv preprint arXiv:2308.06571*, 2023.
- [37] Peiqi Wang, Dongsheng Wang, Yu Ji, Xinfeng Xie, Haoxuan Song, XuXin Liu, Yongqiang Lyu, and Yuan Xie. Qgan: Quantized generative adversarial networks. *arXiv preprint arXiv:1901.08263*, 2019.
- [38] Wei Wen, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. Learning structured sparsity in deep neural networks. In *Advances in neural information processing systems (NeurIPS)*, pages 2074–2082, 2016.
- [39] Yushu Wu, Yifan Gong, Pu Zhao, Yanyu Li, Zheng Zhan, Wei Niu, Hao Tang, Minghai Qin, Bin Ren, and Yanzhi Wang. Compiler-aware neural architecture search for on-mobile real-time super-resolution. In *European Conference on Computer Vision*, pages 92–111. Springer, 2022.
- [40] Zhen Xing, Qi Dai, Han Hu, Zuxuan Wu, and Yu-Gang Jiang. Simda: Simple diffusion adapter for efficient video generation. *arXiv preprint arXiv:2308.09710*, 2023.
- [41] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large video description dataset for bridging video and language. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5288–5296, 2016.
- [42] Changdi Yang, Pu Zhao, Yanyu Li, Wei Niu, Jiexiong Guan, Hao Tang, Minghai Qin, Bin Ren, Xue Lin, and Yanzhi Wang. Pruning parameterization with bi-level optimization for efficient semantic segmentation on the edge. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15402–15412, 2023.
- [43] Zheng Zhan, Yifan Gong, Pu Zhao, Geng Yuan, Wei Niu, Yushu Wu, Tianyun Zhang, Malith Jayaweera, David Kaeli, Bin Ren, et al. Achieving on-mobile real-time super-resolution with neural architecture and pruning search. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4821–4831, 2021.
- [44] Zheng Zhan, Zhenglun Kong, Yifan Gong, Yushu Wu, Zichong Meng, Hangyu Zheng, Xuan Shen, Stratis Ioannidis, Wei Niu, Pu Zhao, and Yanzhi Wang. Exploring token pruning in vision state space models. *arXiv preprint arXiv:2409.18962*, 2024.
- [45] David Junhao Zhang, Jay Zhangjie Wu, Jia-Wei Liu, Rui Zhao, Lingmin Ran, Yuchao Gu, Difei Gao, and Mike Zheng Shou. Show-1: Marrying pixel and latent diffusion models for text-to-video generation. *arXiv preprint arXiv:2309.15818*, 2023.
- [46] Shiwei Zhang, Jiayu Wang, Yingya Zhang, Kang Zhao, Hangjie Yuan, Zhiwu Qin, Xiang Wang, Deli Zhao, and Jingren Zhou. I2vgen-xl: High-quality image-to-video synthesis via cascaded diffusion models. *arXiv preprint arXiv:2311.04145*, 2023.
- [47] Yihua Zhang, Yuguang Yao, Parikshit Ram, Pu Zhao, Tianlong Chen, Mingyi Hong, Yanzhi Wang, and Sijia Liu. Advancing model pruning via bi-level optimization. *Advances in Neural Information Processing Systems*, 35:18309–18326, 2022.
- [48] Pu Zhao, Fei Sun, Xuan Shen, Pinrui Yu, Zhenglun Kong, Yanzhi Wang, and Xue Lin. Pruning foundation models for high accuracy without retraining. *arXiv preprint arXiv:2410.15567*, 2024.
- [49] Daquan Zhou, Weimin Wang, Hanshu Yan, Weiwei Lv, Yizhe Zhu, and Jiashi Feng. Magicvideo: Efficient video generation with latent diffusion models. *arXiv preprint arXiv:2211.11018*, 2022.

# Appendix

## A Memory Snapshot during inference

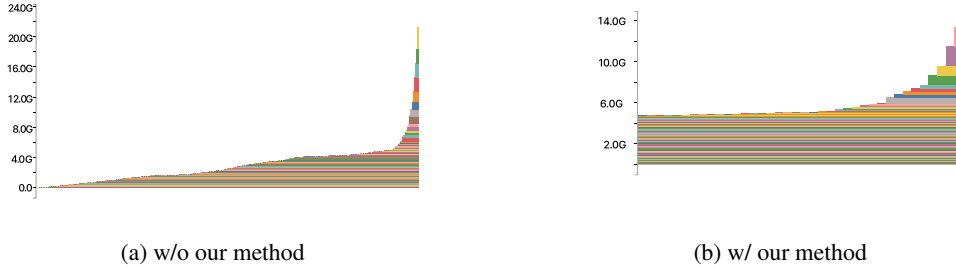


Figure A1: GPU memory snapshot of active cached segment timeline for Stable Video Diffusion with 14 frames  $512 \times 512$

We provide memory snapshots under different configurations during inference, demonstrating the effectiveness of memory reduction. An example is shown in Fig.A1. This example shows the memory reduction of our method on SVD with  $512 \times 512$  resolution. The snapshot is collected following the tutorial<sup>2</sup>.

## B Key Step Search for Step Rehash

**Example of Step Rehash.** For step  $s + 1$ , we only conduct part of the computations in the final up-sampling block, skipping most of the computations in the U-net. Similarly, we can skip multiple steps. For example, if we skip both step  $s + 1$  and  $s + 2$ , to obtain the output  $U_3^{s+2}$  for step  $s + 2$ , we feed the output features  $U_3^{s+1}$  into the final up-sampling block of step  $s + 2$ , where  $U_3^{s+1}$  is also obtained from  $U_3^s$  following the above reusing and skipping method.

**Similarity patterns.** The feature similarity between different steps exhibits certain patterns. As shown in Fig. 5a, at initial steps, the similarity is high (above 97%) across multiple steps such as from step 0 to step 13. In the middle steps, the high similarity only appears within a small step range. For example, the similarity between step 17 and step 19 is lower than 93%. In the final steps, the high similarity appears in a slightly larger step range, such as from step 20 to step 22, with above 93% similarity.

We address the where-to-skip problem with a fixed strategy to skip the computations from the specific blocks. Next, we address the when skip problem to choose what steps can be skipped based on the similarity map. Given the similarity map as shown in Fig. 5a, the similarity value between step  $i$  and  $j$  can be represented by  $S_{ij}$  as shown in the similarity map. We develop a search method to find the key steps with feature rehash and skip the other steps.

The algorithm is shown in Algorithm 1. We use a threshold to select the key steps. If the similarity of multiple consecutive steps is above the threshold, we only select the start and end steps as key steps, and the middle steps can be skipped. Typically, a larger threshold leads to more key steps with high generation performance close to the original one, and a smaller threshold leads to skipping more steps and, thus, computations with faster generation.

---

### Algorithm 1 Key step search in step rehash

---

**Require:** The similarity map  $\mathbf{S}$ , the similarity threshold  $\gamma$ , the maximum step number  $K$

**Ensure:** The set of key steps  $G$

$i \leftarrow 0, j \leftarrow 0, G \leftarrow \{\}, G \leftarrow G \parallel i$

**while**  $i < K$  **do**

**if**  $S_{ij} \geq \gamma$  **then**

$i \leftarrow i + 1$

**else**

$G \leftarrow G \parallel i$

$j \leftarrow i$

$G \leftarrow G \parallel K - 1$

**return**  $G$

---

<sup>2</sup><https://pytorch.org/blog/understanding-gpu-memory-1/>

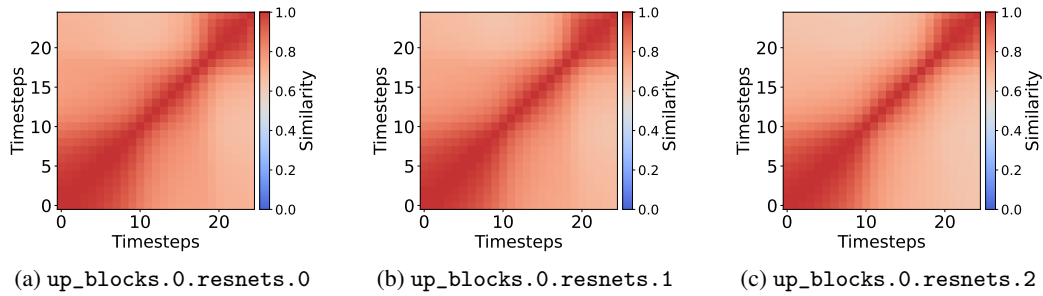


Figure A2: Similarity maps of different temporal layers in `up_blocks.0.resnets`.

We provide sample PyTorch snippets for operation grouping and Step Refresh. The sample code effectively reduces the peak memory and accelerates the inference speed. However, the pipeline is not released because it requires specific compilation support.

### C Similarity map of middle layers

We illustrate the similarity map of several layers closer to the *mid-block* of the UNet, showing that the similarity of these layers is relatively low compared to the results in Fig. 5.

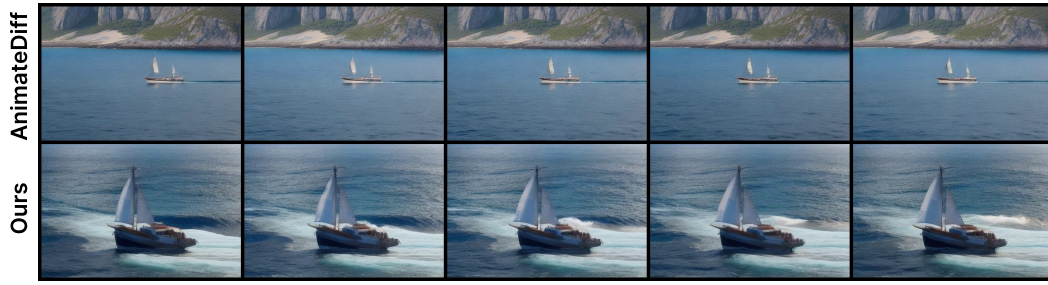
### D Visual comparison with DeepCache

We provide visual comparison of our method with DeepCache in here. As we can see, our method produces more vivid and detailed sample than DeepCache.



Figure A3: Visual comparison of our method with DeepCache.

## E More quality results



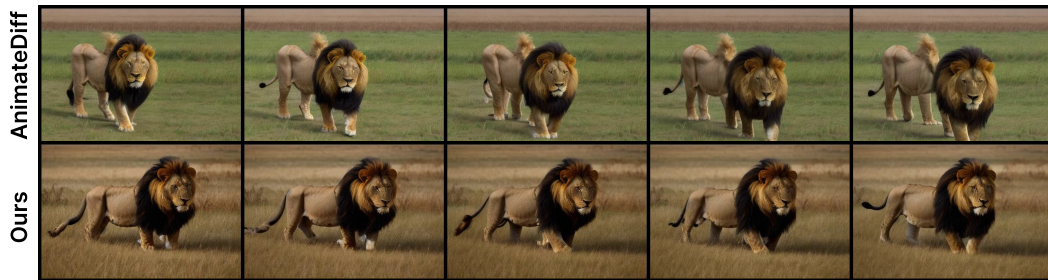
*a boat sailing in the ocean with a rocky landmass in view*



*a bus pulls up to a curb then pulls off*



*the expensive audi car going very fastly in the road at the center of the forest*



*a lion is shown with another lion walking through the field*

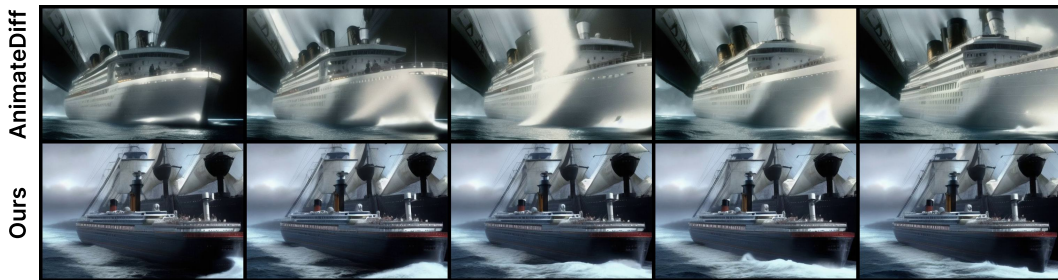
Figure A4: Quality evaluation of using our method on baseline models.



*a bunch of loud lamborghini s a driving up and down the street all at different times*



*2 cute parrots sitting nicely*



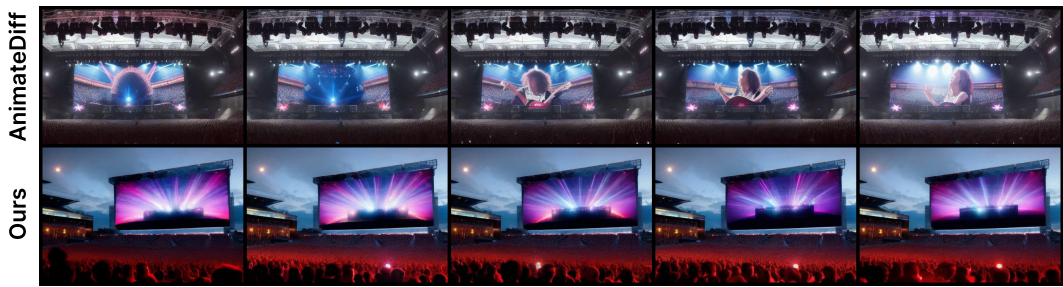
*a movie trailer of james cameron s titanic movie*



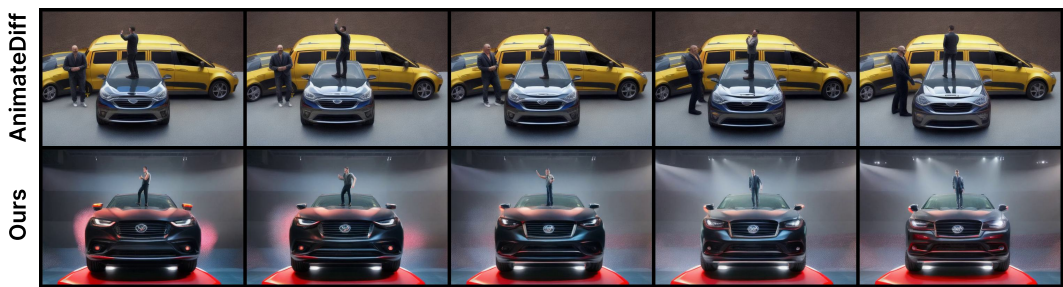
*a street with traffic shown first then a setting of custom made convertables are being displayed*

Figure A5: Quality evaluation of using our method on baseline models.





*a band is performing and being shown on a large screen*



*a man is standing on a car at a car tv show*



*a campaign ad for hillary clinton*



*a beautiful rhyme about the roadbus journey for children*

Figure A6: Quality evaluation of using our method on baseline models.

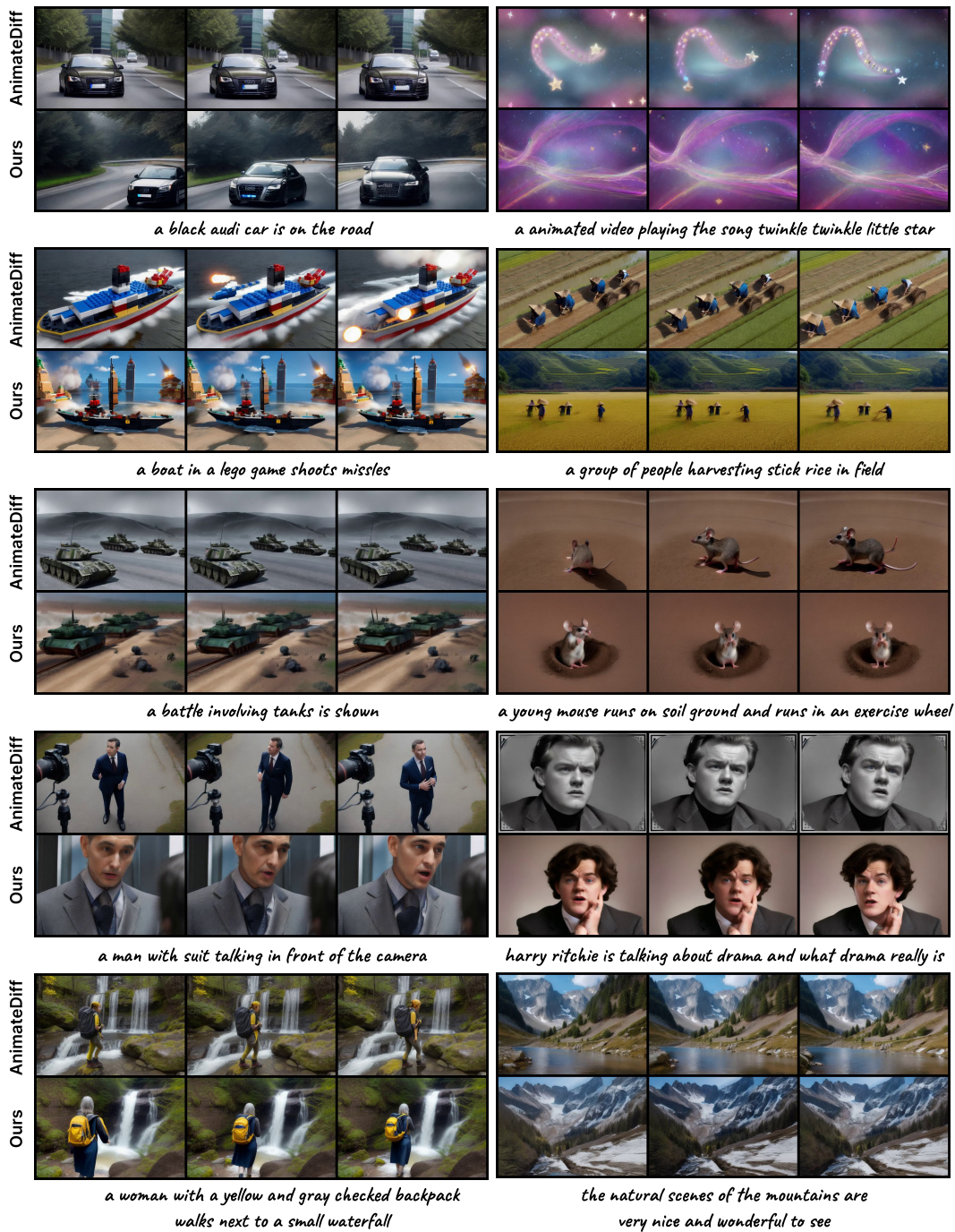


Figure A7: Quality evaluation of using our method on baseline models.

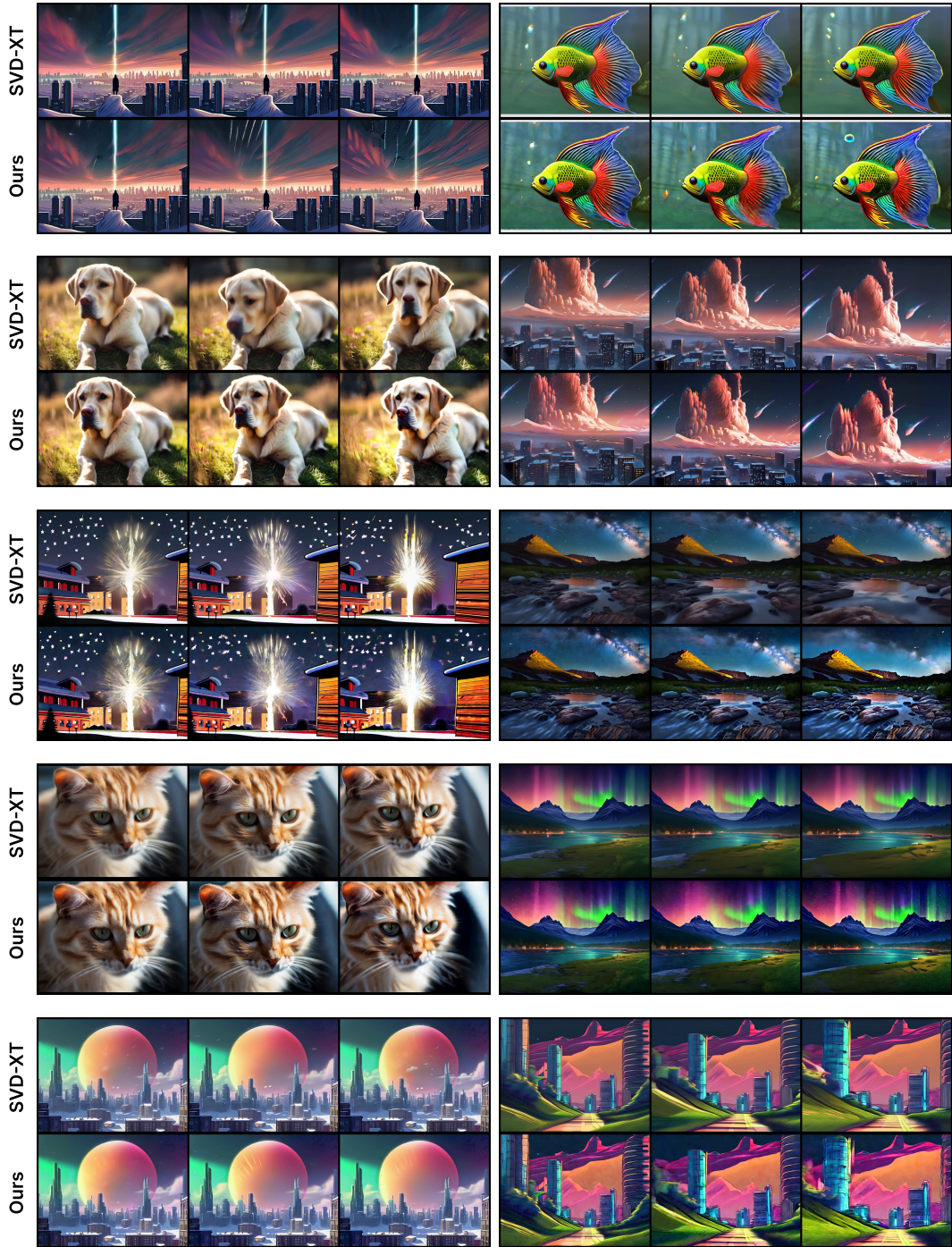


Figure A8: Quality evaluation of using our method on baseline models.

## NeurIPS Paper Checklist

### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: We explain method and summarize the contribution in introduction.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: The limitation is included in conclusion section.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

### 3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: Our paper does not include theoretical results.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: This paper fully discloses all the information needed to reproduce the main experimental results.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

#### 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The datasets and models we used is open-source, and we will open source our code soon.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We have specified all the training and test details necessary to understand the results

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

## 7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We report average results of multiple runs in our experimental section. Our paper does not report error bars.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)

- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We explain the computation resources in experiment section.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

## 9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

Answer: [Yes]

Justification: Research is conducted in the paper conform with NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

## 10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: Our paper is not highly related to societal impacts.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: Our paper poses no such risks. Our work does not release a new model.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: CC-BY 4.0, and we referenced the works that we used to implement our code.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.



- If this information is not available online, the authors are encouraged to reach out to the asset’s creators.

### 13. **New Assets**

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: Our paper does not release new assets.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

### 14. **Crowdsourcing and Research with Human Subjects**

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: This paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

### 15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: This paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.