

4DVD: Cascaded Dense-view Video Diffusion Model for High-quality 4D Content Generation

Shuzhou Yang¹, Xiaodong Cun², Xiaoyu Li^{3*}, Yaowei Li¹ and Jian Zhang^{1*}

¹Peking University Shenzhen Graduate School, Shenzhen, China.

²Great Bay University, Dongguan, China.

³Tencent, Shenzhen, China.

*Corresponding author(s). E-mail(s): xliea@connect.ust.hk; zhangjian.sz@pku.edu.cn;
Contributing author(s): szyang@stu.pku.edu.cn; cun@gbu.edu.cn; liyaowei01@gmail.com;

Abstract

Given the high complexity of directly generating high-dimensional data such as 4D, we present 4DVD, a cascaded video diffusion model that generates 4D content in a decoupled manner. Unlike previous multi-view video methods that directly model 3D space and temporal features simultaneously with stacked cross view/temporal attention modules, 4DVD decouples this into two subtasks: coarse multi-view layout generation and structure-aware conditional generation, and effectively unifies them. Specifically, given a monocular video, 4DVD first predicts the dense view content of its layout with superior cross-view and temporal consistency. Based on the produced layout priors, a structure-aware spatio-temporal generation branch is developed, combining these coarse structural priors with the exquisite appearance content of input monocular video to generate final high-quality dense-view videos. Benefit from this, explicit 4D representation (such as 4D Gaussian) can be optimized accurately, enabling wider practical application. To train 4DVD, we collect a dynamic 3D object dataset, called D-Objaverse, from the Objaverse benchmark and render 16 videos with 21 frames for each object. Extensive experiments demonstrate our state-of-the-art performance on both novel view synthesis and 4D generation. Our project page is <https://4dvd.github.io/>

Keywords: 4D generation, multi-view diffusion model, video generation

1 Introduction

Dynamic 3D object generation (4D generation) aims to simultaneously create realistic 3D object content and its motion in 3D space. Considering that the 3D world we live in is inherently dynamic (such as fluttering flags, walking people, and moving objects), 4D generation is crucial to achieving immersive content creation and visual content experience. In this work, we attempt to generate 4D objects from widely available monocular videos, which enables effortless and advanced AR/VR creation.

Generating 4D asset from a given monocular video is challenging as this means the model needs to simultaneously reason about the appearance and motion around the object in 3D space based on a single 2D view, which is an ill-posed problem. To this end, recent attempts have been made to reconstruct 4D content through Score Distillation Sampling (SDS) strategy [1–5], which distill the 3D priors from multi-view or image generation models and temporal priors from video diffusion models. However, these methods require a time-consuming distillation process and tend to output over-smooth results with the Janus problem. Furthermore, the priors these methods distill

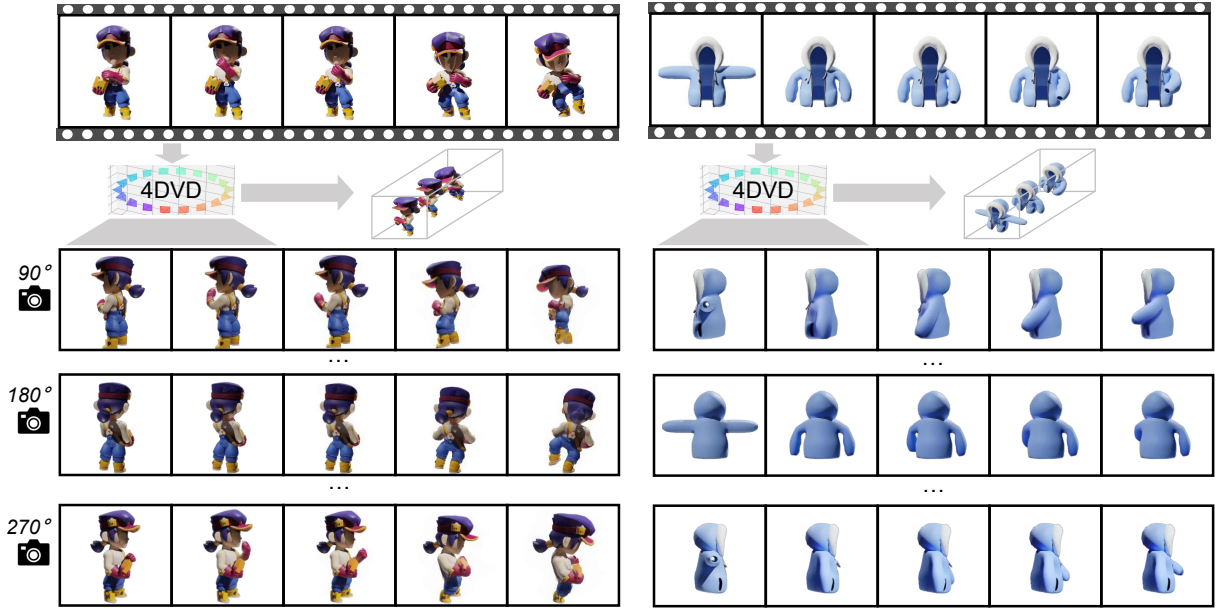


Fig. 1: 4DVD takes a monocular video as input and generates multi-view videos. Benefiting from unprecedentedly dense-view modeling, the generated results maintain high spatial and temporal consistency.

come from different models, *i.e.*, videos from video generation models, and views from image generation models, which may conflict with each other during the optimization. Another line of works [6–8] instead unify the prediction of multi-view and motion in a single model. Most commonly, they propose to first get the multi-view (V) conditions based on the first frame, then generate 4D content combined with the input monocular video (T). The final 4D content can be regarded as an image grid with a size of $T \times V$, and the model itself aims to complete the grid content based on edge conditions. Existing work is implemented through training on massive data, however, directly learning from such high-dimensional data remains difficult. On the one hand, simultaneously learning both 3D space and temporal motion with an end-to-end black-box model increases the difficulty of training. On the other hand, limited VRAM restricts the accuracy of 4D modeling. Since the model has to process $T \times V$ images in every single batch, even with the most advanced GPU, the viewpoint value V can only be set to a small value (such as 4 in [9] and 8 in [6]). The sparse viewpoints inevitably impact the capacity of modeling dynamic 3D space.

To address this problem, we propose **4DVD**, a cascaded dense-view video diffusion model that aims to generate higher quality 4D assets in a decoupled manner. We find that in our case, *i.e.*, multi-view video generation, spatial information is actually over-redundant, since the same visual content oftentimes appears on different views and frames. Hence, 4DVD firstly predicts consistent dense-view layouts at the expense of spatial details, and achieves structure-aware conditional generation based on these, obtaining high-quality 4D assets. Specifically, In the first stage, we downsample multi-view videos to low-resolution versions, enabling training with denser views for more efficient 3D spatial modeling. This trade-off between viewpoint number and resolution brings more prolific priors of 3D content and motion, but remains ambiguous appearance. Therefore, in the second stage, 4DVD takes multi-view layouts as the hint and operates structural conditional generation, producing high-quality multi-view results. To realize this insight, in the first stage, we inherit the structure of SV4D [6], a State-Of-The-Art (SOTA) multi-view video generation model, but extend it to denser view prediction by training on the low-resolution dense-view video data. To achieve structural-aware conditional generation of

the second stage, we carefully design a unique control branch that injects multi-view layout priors to the mainstream layer-by-layer. Considering that the control branch only has multi-view layout information but lacks appearance guidance, we propose Monocular Appearance Propagation (MAP) module to further incorporate the high-quality visual content of the input monocular video in the final 4D conditional generation process. Finally, 4DVD unifies the structure modeling and structure conditional spatio-temporal generation, brings significant performance improvements and can be used for explicit 4D reconstruction. The experimental results prove the effectiveness of our proposed method. We found that open-source 4D datasets contain numerous low-quality cases, such as broken components, minor or over-excessive motion, *etc.* To train 4DVD effectively, we collected a high-quality dynamic multi-view dataset from Objaverse [10], called D-Objaverse. It will be released when this paper is published.

We summarize our key contributions as follows:

- We present 4DVD, a cascaded 4D content generation pipeline that improves generation quality by modeling 3D space and motion from unprecedented dense views.
- A structure-aware spatio-temporal generation branch is proposed to organically combine the input monocular video and dense-view coarse results, injecting these guided conditions to the mainstream model layer-by-layer and efficiently producing high-quality 4D results.
- For effective training, we carefully filter and render 16-view video data from Objaverse, called D-Objaverse. This subset contains complete object structures and natural motions.
- Both quantitative and qualitative experiments prove that 4DVD achieves state-of-the-art performance in both multiview video synthesis and 4D reconstruction.

2 Related Work

2.1 3D Generation from Text or Image

We first discuss 3D generation methods that produce 3D assets from text or images. DreamFusion [11] proposes a Score Distillation Sampling (SDS) method to optimize the 3D content through the distillation of 2D diffusion priors, which becomes a paradigm for 3D generation using 2D diffusion models. However, due to its limited quality and time-consuming generation process, many subsequent approaches [12–22] have been proposed to develop more effective and efficient distillation methods, aiming to address the Janus problem, alleviate over-smoothed content, and speed up the generation process. To overcome the time-consuming optimization process of these SDS-based methods for each generated 3D object, researchers have explored direct prediction of 3D models in a feedforward manner via a large reconstruction model [23–29], which could generate the 3D model instantly. However, due to the high complexity of the 3D representations such as NeRF [30], 3D Gaussians [31], and triplane and limited network capacity, these methods usually could only generate objects with simple geometry and texture. Other works propose multi-view generation models [32–35] that can produce view-consistent multi-view images of the 3D content. Compared with direct 3D generation via a large reconstruction model, reconstructing 3D content using the generated multi-view images greatly reduces the complexity of directly learning 3D representations and achieves superior performance. Inspired by this strategy, we develop a new algorithm that produces consistent multi-view videos (instead of images) to reconstruct 4D objects.

2.2 4D Content Generation

The development of object-centered 4D content generation is similar to that of 3D generation. Optimization-based methods are also firstly explored [1–5, 36–40], using SDS loss to optimize the 4D representation [41–43] through pre-trained video diffusion models. The key insight of these methods is the distillation of a text-to-video diffusion model and a multi-view diffusion model to get

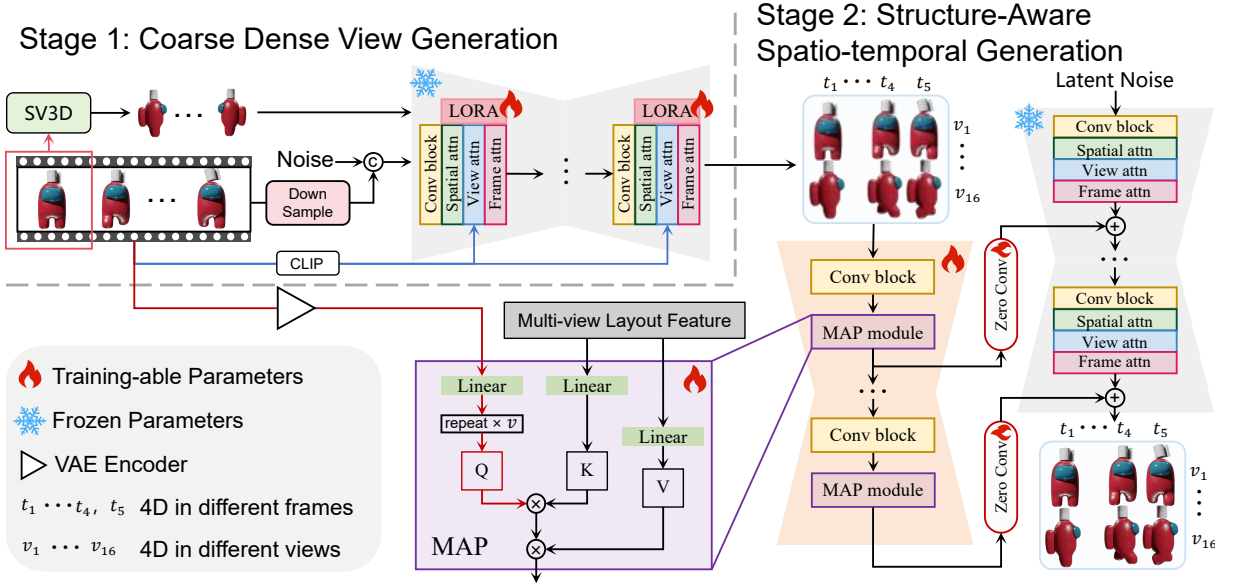


Fig. 2: Pipeline of 4DVD. Our model consists of two stages: The coarse dense view generation receives a monocular video and reference views, producing the coarse 16-view layouts. Based on layout condition, the structure-aware spatio-temporal generation predicts detailed multi-view videos. To get appearance guidance, we take input monocular video as an external condition and develop Monocular Appearance Propagation (MAP) module to integrate it, as shown in the purple region.

temporal and cross-view priors. For example, 4Dfy [3] divides the optimization process into three stages, training based on the multi-view model, image model, and video model, which improves 3D consistency, generation quality, and motion quality respectively. However, these methods can take hours to produce 4D content due to the time-consuming distillation process. In addition, large reconstruction models [7] and multi-view video generation models [6, 9, 44] are also developed. The former directly predicts 3D Gaussian parameters for each frame, which requires enormous GPU resources for training, and its generation quality is limited by complex modeling of implicit prediction parameters. The latter predicts 3D consistency multi-view videos, which require fewer computing resources and enable high-quality generation. However, due to the limited VRAM, current multi-view video generation methods can only model dynamic 3D space with relatively sparse views. This not only restricts the ability of generation models to perceive the full 3D space but also cannot be used to reconstruct explicit 4D representations such as 4D Gaussian accurately, which usually requires dense views as input.

More recently, another line of works [8, 45, 46] attempted to generate 4D scenes containing multiple objects and wide background. But these methods mainly focus on extending monocular scene videos to flexible forward views for immersive watching experience, instead of facilitating dynamic 3D digital asset creation. To produce full 360-degree views for object-centered creation, we propose a cascaded dense-view video generation model to produce highly cross-view consistent multi-view videos with high-quality details.

3 Method

Given a monocular video $\mathbf{I} \in \mathbb{R}^{T \times D}$ of a dynamic object with T frames and $D = 3 \times H \times W$ dimensions, 4DVD aims to faithfully predict the dynamic content in different views with both view and temporal consistency. Specifically, our goal is to extend the video \mathbf{I} from $\mathbb{R}^{T \times D}$ to $\mathbb{R}^{T \times V \times D}$, where V is the number of viewpoints. Considering more viewpoints facilitate more complete modeling for dynamic 3D spaces, in this work, we introduce a coarse-to-fine generation pipeline that

produces highly consistent dense-view videos in a decoupled-cascaded manner.

3.1 Preliminaries: Latent Diffusion Models

Existing multi-view video diffusion models can be regarded as a kind of Latent Diffusion Model (LDM). These models [47, 48] generate visual content in the latent space that can be encoded and decoded by the Variational Auto-Encoder (VAE) model. In training, given the latent \mathbf{z}_0 encoded by the VAE encoder, we add noise ϵ of various scales to it following a predefined schedule defined as:

$$\mathbf{z}_t = \sqrt{\bar{\alpha}_t} \mathbf{z}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, \quad (1)$$

LDM is trained to precisely predict the added noise ϵ in the noisy latent \mathbf{z}_t with a neural network $\epsilon_\theta(\cdot)$, whose objective function can be summarized as:

$$\min_{\theta} \mathbb{E}_{\mathbf{z}_0, \epsilon \sim \mathcal{N}(0, I), t} \|\epsilon - \epsilon_\theta(\mathbf{z}_t, t, \mathbf{c})\|_2^2, \quad (2)$$

where t is the sampling step and \mathbf{c} is the text condition guiding the generated content. After training, $\epsilon_\theta(\cdot)$ receives a random latent noise \mathbf{z}_T , denoises step by step with a prompt condition \mathbf{c} , and decodes the final latent through a VAE decoder to get final results.

3.2 Cascaded Coarse-to-fine Pipeline

Simultaneously generating $T \times V$ images requires massive VRAM and is inefficient. Previous methods that process the $T \times V$ image grid usually restrict to limited viewpoints, which makes it challenging to fully comprehend the geometry and motion in 3D space and reconstruct the 4D representation accurately, as this requires highly consistent dense-view input. To fully unleash the potentiality of generation models to create dynamic content in 3D space, we propose a two-stage cascaded pipeline as shown in Fig. 2 to generate prolific multi-view videos from a monocular input video. Our key idea is to learn dense-view layout priors to model the primitive structure of 4D objects precisely in the coarse stage. And then, guided by these coarse structural priors, we synthesize high-quality final results while maintaining the view

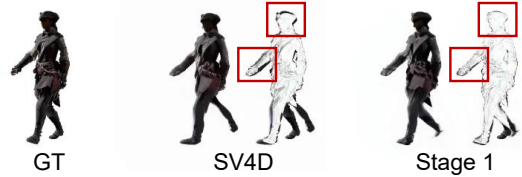


Fig. 3: Visual and residual results. We provide the GT image, results of SV4D and our first stage model, and the corresponding residual images compared to the GT image. Enforcing SV4D to produce 16 views leads to large errors (as shown in regions boxed in the red), while results from our first stage achieve better results after fine-tuning.

and temporal consistency in the fine stage. Note that since spatial information is redundant across view and time, resolution can actually be sacrificed to learn multi-view layout priors, allowing for more adequate and accurate 3D modeling from more viewpoints. These dense views bring more consistent generation results for 3D content and object motion. In the fine stage, we dynamically select a subset of the dense views in each training iteration to supervise under limited VRAM. Considering that the fine stage aims to perform spatio-temporal generation guided by the layout hints from the coarse stage, instead of purely generate the dynamic 3D space from scratch, it is sufficient and effective to optimize with only a part of viewpoints in each iteration. During inference, due to the absence of back propagation, 4DVD can output the complete multi-view videos based on the full 16-view layout guidance in the fine stage. This coarse-to-fine process aligns with the workflow of the human visual system: we first perceive the general layout of a dynamic object at a glance and then focus on enriching its details.

3.3 Coarse Dense View Generation

Previous multi-view video generation methods like SV4D [6, 49] are trained on sparse view videos, which could also be adapted to predict more views by directly adjusting their hyperparameters during inference. But this naive modification has some problems. First, it fails to predict accurate results when the number of viewpoints is greater than the number of views in the training set. As illustrated in Fig. 3, we present visual results produced

by SV4D with 16 viewpoints and their Ground Truth (GT) and visualize their residual image on the right. One can see that although SV4D could generate more views, its output exhibits abnormal geometry such as the distorted right leg, and significantly deviates from the GT image. Highlighted in red boxes of Fig. 3, the contours of the residual image exhibit obvious differences, indicating that the overall output layout of SV4D does not accurately match the GT view. This is because SV4D is restricted by the sparse training views (8 views on the most advanced GPU), making it challenging to fully model the 4D space. Forcing generalization and producing more views exposes this inability to accurately depict 4D content.

To address this problem, we propose to generate denser view videos at the expense of spatial content. We finetune SV4D on dense view low-resolution videos (e.g., 16 views in our setting) using LORA [50] to fit the VRAM requirements. Our goal is not to directly predict the 16-views spatio-temporal information, but only to obtain the 4D layout priors that are consistent across views. As shown in Fig. 2, given a monocular video as input, we initially employ the existing image-to-3d method [51] to synthesize novel views for the first frame. Taking the reference video and 16 reference views as input, the network in the coarse stage predicts the $T \times V$ image grid at a low resolution. Considering spatial messages are actually over-redundant in our case, we believe that sacrificing it in exchange for a more adequate modeling of 4D space is worthwhile and effective. We incorporate the LORA module into three attention blocks, which are used to process spatial features, fuse with the input video, and fuse with the reference views, respectively.

As shown in Fig. 3, the output of Stage 1 exhibits reliable geometry, and its residual image exhibits less difference, demonstrating that the unique design in our coarse stage leads to more accurate 4D modeling and generation. Theoretically, we could use even more viewpoints for training by further reducing video resolution. However, we observe that the VAE encoder-decoder [47] cannot reproduce images with too low resolution faithfully. We believe this is because it is pre-trained on high-quality images with high resolutions. Consider that 16 views already cover 4D object densely enough, and further reducing spatial messages has the potential risk of even

erasing desired layout information. According to our experiments, we choose 256×256 as the final resolution for our first stage, which is enough to enable dense-view training.

3.4 Structure-Aware Spatio-temporal Generation

Since we produce multi-view videos simultaneously with attention mechanisms, they preserve consistent structure information spatially and temporally. However, these layouts exhibit blurred textures, low resolution, and quality that deviate from the input video. Therefore, our goal in the fine stage is to produce high-quality dense view videos based on these consistent layouts, and we introduce structure-aware spatio-temporal generation to achieve this.

Inspired by previous conditional generation methods [52–54], which extract layout features through a specific branch and inject them to the base model, we also utilize SV4D as mainstream model and introduce a trainable structure condition branch. It encodes the coarse multi-view layouts to latent space that is aligned to mainstream model, and injects them into the base model. As shown in Fig. 2, we take the output of the first stage as the input condition of our structure condition branch, which is stacked with convolution blocks and our proposed MAP modules. During training, we freeze the base model and update the parameters of the structure condition branch only. Constrained by the limited VRAM, we randomly select 4 views of the 576×576 resolution at each training iteration. Since the condition branch injects the aligned layout features into the base model, and the goal of the second stage is not to model 4D structure but to provide better layout guidance for each view, there is no need to provide full-view condition.

Note that this stage aims to generate multi-view videos with elegant appearance, so high-quality appearance prior is required but so far only multi-view layout conditions are available. Therefore, we utilize the input monocular video as external condition and propose Monocular Appearance Propagation (MAP) module to propagate its appearance to multiple views. As shown in Fig. 2, in each MAP module, the high-quality input video is injected to provide appearance reference. The details of the MAP module are

given in the purple region, which is a basic cross-attention design with some unique modifications. Note that conventional cross-attention uses the features of the current model to calculate the \mathbf{Q} matrix and employs external conditions for \mathbf{KV} matrices. But in our case, we calculate \mathbf{KV} matrices with the multi-view features and use reference video for \mathbf{Q} matrix. This is because our condition branch features structure information of multi-view videos whose shape is $T \times V$, but the reference video only records the high-quality content of a single view ($T \times 1$). Therefore, MAP module actually aims to integrate these two unaligned features. Focusing on this special issue, we swap the calculation of \mathbf{Q} and \mathbf{KV} , enabling the model to extend high-quality reference video to novel views based on the condition of multi-view layouts, which is easier to train than conventional operations that forcibly upscale the multi-view coarse results based on monocular video. Experiments in Sec. 5.5 prove the effectiveness of our unique design.

To enable this modified cross-view calculation, we repeat the reference video for V times before multiplication. Denote the linear layers that calculate \mathbf{Q} , \mathbf{K} , \mathbf{V} as L_Q , L_K , L_V respectively, the operation of MAP can be expressed as:

$$\begin{aligned}\mathbf{Q} &= L_Q(F^{ref}) \in \mathbb{R}^{WHV \times N_c}, \\ \mathbf{K} &= L_K(F^{LR}) \in \mathbb{R}^{N_c \times WHV}, \\ \mathbf{V} &= L_V(F^{LR}) \in \mathbb{R}^{WHV \times N_c},\end{aligned}\quad (3)$$

where WH and N_c mean the spatial resolution and feature channel number of the current layer, F^{ref} is the normalized VAE feature of the reference video, and F^{LR} means the coarse multi-view layout features. For training, we set $V = 4$ and during inference, we use the full coarse image grid as the condition, *i.e.*, $V = 16$.

3.5 Training

We train our cascaded model in a decoupled manner for efficiency. In the first stage, we enforce the model to produce 16-view videos with the conditions of reference video \mathbf{c}_f and view \mathbf{c}_v , objective function can be expressed as:

$$\min_{\theta} \mathbb{E}_{\mathbf{z}_0, \varepsilon \sim \mathcal{N}(0, I), t} \|\varepsilon - \varepsilon_{\theta}(\mathbf{z}_t, t, \mathbf{c}_f, \mathbf{c}_v)\|_2^2, \quad (4)$$

where θ is the parameters of our added LORA. In the second stage, we optimize the condition branch to realize structure-aware spatio-temporal generation, which is conditioned on the randomly sampled 4 coarse layouts $\mathbf{I} \in \mathbb{R}^{T \times 4}$. Using θ to represent trainable parameters and loss function is formulated as:

$$\min_{\theta} \mathbb{E}_{\mathbf{z}_0, \varepsilon \sim \mathcal{N}(0, I), t} \|\varepsilon - \varepsilon_{\theta}(\mathbf{z}_t, t, \mathbf{I})\|_2^2, \quad (5)$$

4 D-Objaverse Data

Although Objaverse contains numerous 3D objects, most of them are static and cannot be used to train multi-view video models. Some open-source dynamic 3D data is available such as ObjaverseDy used by SV4D or Consistent4D. However, we notice that these available dynamic 3D data either only contain limited examples or have quite a lot of low-quality cases. These low qualities can be summarized into the following three manifestations: **1)** Incomplete structure. Many 3D objects are just parts of common objects, such as an arm or a leg. **2)** Abnormal motion. Some cases exhibit excessive or unnoticeable movement. **3)** Undefined content. Some examples show blurry/chaotic content, although they are also dynamic 3D objects. To effectively train multi-view video model, we filter out high-quality dynamic 3D assets and render high-resolution videos in 16 viewpoints. We have designed three filtering indicators specifically for the three main low-quality situations. **1)** We use LPIPS to select samples of acceptable quality. **2)** We filter out samples with reasonable motion by calculating the optical flow between frames. **3)** And at last, we utilize a multimodal LLM (*i.e.*, CogVLM [55]) to filter out examples containing incomplete structures. We collected 41k dynamic 3D assets based on ObjaverseDy, containing 27k samples, and other accessible dynamic 3D data, and finally filtered out 17k samples with 16-view videos that met our criteria. This high-quality dynamic 3D dataset is called D-Objaverse since it is mainly filtered from Objaverse. As a major contribution of 4DVD, D-Objaverse will be released when 4DVD is published. We showcase some frames of front-view videos from D-objaverse in Fig. 4.



Fig. 4: Some cases in D-Objaverse. We carefully filtered and rendered a high-quality dynamic 3D datasets from Objaverse data, called D-Objaverse. Here we showcase some front-view cases of D-Objaverse.

5 Experiments

5.1 Implementation Details

Hyperparameter Settings.

We train and evaluate the performance of different methods on the collected D-Objaverse dataset. For a fair comparison, we randomly select 50 dynamic objects for evaluation and the remaining ones are used for training. We use Blender to render 16-view videos with the resolution of 576×576 for each 4D asset. Each video clip is 10 FPS and is downsampled to a corresponding low-resolution version (*i.e.*, 256×256) for the coarse stage. During training, we set $T = 5$. For the first stage, we supervise the network with 16-view low-resolution videos. For the second stage, we randomly select 4 views from the high-resolution videos and make the model to produce a 5×4 high-resolution image grid in each iteration, where the corresponding 4 views of low-resolution version are adopted as structural generation conditions. During inference, the model can predict the entire high-resolution 5×16 image grid at once within 30GB VRAM. To extend to 21-frames full videos, we follow the same anchor-sampling approach as SV4D. D-Objaverse will be released when this paper is published.

Metrics.

Following previous methods [6, 7], we evaluate the generation quality of different methods with three

widely used metrics: Learned Perceptual Similarity [56] (*LPIPS*), CLIP-score (*CLIP-S*), and *FVD*. Considering the characteristics of $T \times V$ 4D generation, SV4D [6] proposes some variants of *FVD* and we adopt three of them, 1) *FVD-F*: fix the viewpoint and calculate *FVD* across frames. 2) *FVD-V*: fix the frame and calculate *FVD* across views. 3) *FVD-Diag*: calculate *FVD* over the diagonal images of the $T \times V$ image grid. These three metrics could evaluate the 4D content in terms of time, space, and space-time, respectively.

Baselines.

We compare our method with other four well-known related methods for 4D generation from a monocular video, including SDS-based method STAG4D [38], multi-view video generation models 4Diffusion [9], SV4D [6], and dynamic large reconstruction model L4GM [7].

5.2 Quantitative Comparison

Multi-view Video Synthesis.

Considering that our method can achieve 16-view generation but others cannot, we only use videos from the four orthogonal viewpoints for comparison. As shown in Tab. 1, we calculate evaluation metrics on the $T(21) \times V(4)$ image grids produced by different methods. For STAG4D and L4GM, we render the same 4 views from their outputs to get image grids. One can see that even when only compared to the orthogonal 4-view videos, 4DVD



Fig. 5: Visual comparison of the generated multi-view videos. We show two timesteps and two views for each example. Compared to other baseline methods, 4DVD enables higher quality generation which has finer textures, reasonable appearance, and better consistency.

Table 1: Quantitative comparison of produced multi-view videos and produced 4D Assets. 4DVD achieves superior performance in visual quality, video frame consistency and multi-view consistency. The best results are highlighted in **bold**.

| Model | Type | Multi-View Videos | | | | | 4D Assets | | |
|------------|----------------------------|-------------------|--------------|---------------|---------------|---------------|--------------|--------------|---------------|
| | | LPIPS ↓ | CLIP-S ↑ | FVD-F ↓ | FVD-V ↓ | FVD-Diag ↓ | LPIPS ↓ | CLIP-S ↑ | FVD ↓ |
| STAG4D | Optimization-based | 0.156 | 0.892 | 714.51 | 528.06 | 652.17 | 0.157 | 0.902 | 879.19 |
| L4GM | Large reconstruction model | 0.144 | 0.910 | 606.90 | 467.81 | 525.22 | 0.143 | 0.911 | 543.57 |
| 4Diffusion | Multi-view video model | 0.155 | 0.899 | 853.89 | 776.92 | 816.88 | 0.165 | 0.874 | 822.29 |
| SV4D | Multi-view video model | 0.136 | 0.913 | 586.41 | 482.64 | 542.87 | 0.161 | 0.897 | 677.56 |
| 4DVD | Multi-view video model | 0.133 | 0.927 | 507.12 | 314.44 | 456.01 | 0.136 | 0.919 | 438.41 |

still achieves better results on both cross-view consistency and visual quality.

4D Generation.

To compare the quality of 4D assets, we use the same 4D Gaussian representation following DreamGaussian4D [5] to reconstruct for multi-view video methods. For STAG4D [38] and L4GM [7], we directly use their outputs. After

getting the 4D Gaussian of different methods, we could render videos with both camera view and temporal index change simultaneous for comparison in 4D. The details of rendered free-view videos can be found in supplementary materials. As we evaluate the free-view video instead of the image grid, we could directly use *FVD* to measure spatio-temporal quality. We report the quantitative comparison between our results and

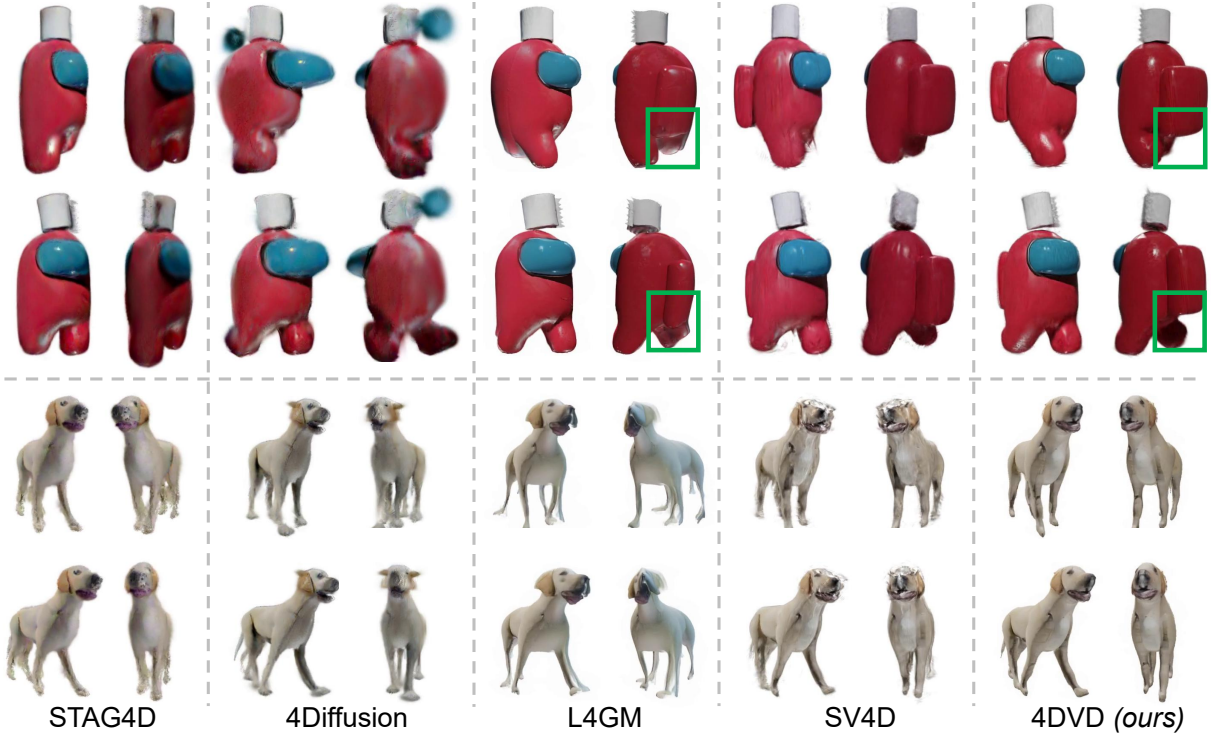


Fig. 6: Visual comparison of the generated 4D outputs. We show two timesteps and two views for each example. By leveraging the dense-view videos generated by our method, 4D Gaussian can be reconstructed well with high-quality geometry and textures. Compared to prior works, our results are more consistent and detailed with a more reasonable 3D structure and satisfactory appearance.

the baselines in Tab. 1. Our method performs best in terms of all metrics, demonstrating its superiority in visual quality (*LPIPS*, *CLIP-S*), 3D and temporal consistency (*FVD*). Although 4Diffusion and SV4D can produce high-quality videos, their limited sparse views lead to visual degradation of their 4D Gaussian results. In contrast, our method enables 16-view generation, which not only brings more consistent results, but also allows high-quality 4D reconstruction. One can see that among all SOTA methods, 4D assets of 4DVD achieves the best scores across all metrics.

5.3 Visual Comparison

Multi-view Video Synthesis.

In Fig. 5, we showcase the video results produced by different methods on three examples with two different angles and timesteps. Notice that our method enables 16-view generation and other methods cannot, for comparison, the selected two viewpoints are orthogonal, which can be directly

produced by all approaches. We observe that the SDS-based methods (*e.g.*, STAG4D) fail to produce reasonable 3D geometries and suffer from the Janus problem. Take the example in the middle of Fig. 5, where we show the results in back and side views, but the back view of STAG4D still has the frontal head, and the side view is only a slice. 4Diffusion, which is a multi-view video model, enables multi-view generation while maintaining consistency across views and frames. However, its results exhibit weird structures and unnatural appearance. The most recent L4GM and SV4D achieve relatively higher quality, but L4GM cannot predict geometry well. We believe this is because directly predicting 3D Gaussian parameters for all frames is complicated, which causes L4GM to be unable to accurately predict the novel view content of cases given in Fig. 5. Meanwhile, SV4D tends to produce over-smooth texture, leading to unsatisfactory appearances. In contrast, benefiting from the proposed decoupled generation strategy, our cascaded model simplifies

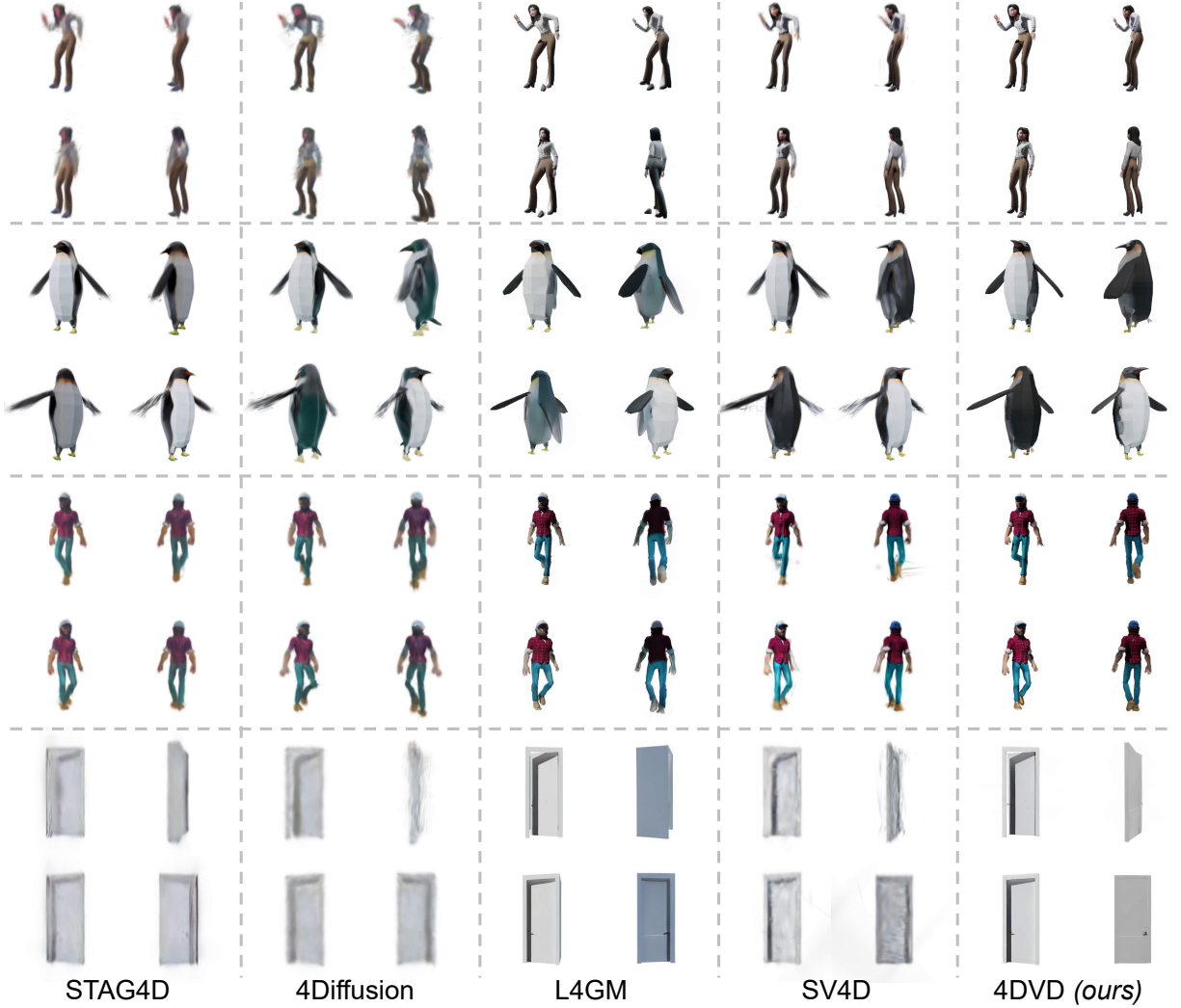


Fig. 7: More visual comparison of the generated 4D outputs from different methods. Similar to Fig. 6, we show two time steps and two views for each example. Compared to prior works, our results exhibit clearer appearance and more prolific geometric structures. We believe this performance advantage mainly comes from the more consistent multi-view video and unprecedented dense viewing angles produced by 4DVD.

the difficulty of modeling 3D content and motion, producing better visual results.

4D Generation.

For the 4D reconstruction results, as shown in Fig. 6, STAG4D still suffers from the Janus problem and fails to model the motion in 3D space well. 4Diffusion mistakes the 3D content and produces distorted structures. For the cases such as the dog with slender legs, L4GM fails to model it since this model generates Gaussian parameters

based on only 4 views, which is sparse for 4D modeling. SV4D only produces 8-view videos, which is not enough to reconstruct the explicit 4D representation. In contrast, our method that directly generates 16 consistent dense view videos enables more accurate and consistent 4D modeling, which can be used to reconstruct 4D Gaussians with high quality. Moreover, to provide stronger evidence, we further display more comparison cases in Fig. 7. One can see that previous work cannot produce high-quality 4D assets while 4DVD

Table 2: Comparison on runtime. The best results are highlighted in **bold**.

| Methods | STAG4D | L4GM | 4Diffusion | SV4D | 4DVD |
|--------------|--------|------|------------|------|--------------------|
| Runtime(s) ↓ | 4200 | 15 | 480 | 565 | 381(51+330) |

enables satisfactory 4D generation. We believe this is mainly due to two reasons. First, due to the advanced novel workflow, 4DVD can produce multi-view videos that are highly consistent in time and viewpoint, which enables accurate 4D reconstruction. Second, 4DVD produces unprecedented dense views, providing ample supervision for 4D modeling.

5.4 Runtime Comparison

Theoretically, as a cascaded diffusion model, 4DVD may require more runtime cost compared with SV4D due to its two-stage coarse-to-fine generation process. However, we find that it is precisely due to the proposed coarse-to-fine generation that 4DVD actually has a runtime advantage over SV4D. We report the runtime cost of different methods on a single A6000 GPU in Tab. 2. One can see that L4GM realizes the fastest speed as it is not a diffusion model and does not require iterative sampling. Meanwhile, 4DVD is even faster than SV4D. The reason is that, SV4D needs to generate from scratch, so its official code chooses to iterate sampling for 20 steps. In contrast, 4DVD generates in a decoupled and cascade manner. In its first stage, we only aim to produce coarse results with low resolution. Lower resolution allows for faster feedforward, and can be generated within only 10 sampling steps. In the second stage, we have obtained the coarse results as the control hint. Since 4DVD benefits from coarse multi-view priors of the first stage, we also set sampling number to 10, which is enough to complete generation. The runtime cost of both stages are listed in the brackets (first + second) of Tab. 2.

5.5 Ablation Study

To validate the effectiveness of our full model, we conduct ablation experiments on the proposed cascaded pipeline. We replace our structure-aware conditional generation branch with existing well-known methods [52, 53], and compare their performance to validate the superiority of our design.

Table 3: Quantitative results of ablation study. The best results are highlighted in **bold**.

| Setting | LPIPS ↓ | CLIP-S ↑ | FVD-F ↓ | FVD-V ↓ | FVD-Diag ↓ |
|----------------|--------------|--------------|---------------|---------------|---------------|
| Stage 1 | 0.147 | 0.901 | 668.25 | 422.56 | 605.73 |
| w/ ControlNet | 0.141 | 0.920 | 558.42 | 407.75 | 535.14 |
| w/ T2I-adapter | 0.173 | 0.896 | 802.40 | 609.45 | 714.37 |
| MAP w/ feature | 0.152 | 0.912 | 572.91 | 397.36 | 512.15 |
| Full Model | 0.133 | 0.927 | 507.12 | 314.44 | 456.01 |

Moreover, we replace the proposed MAP with the common cross-attention mechanism (*i.e.*, “MAP w/ feature”) to demonstrate the effectiveness of our MAP module.

Ablation on condition branch.

We replace our structure condition branch with well-known condition-guided methods such as ControlNet [52] and T2I-adapter [53]. The detailed architectures of these ablation settings can be found in Fig. 8, where we apply our proposed control branch, T2I-adapter, and ControlNet to the mainstream model to realize coarse-to-fine 4D generation in Fig. 8(a), (b), (c), respectively. We compare their visual results in Fig. 9. Stage 1 produces coarse results with low visual quality, which mainly aims to provide multi-view layout priors. To generate spatio-temporal content based on these, we apply various architecture to find a feasible injection strategy. Using a T2I-Adapter that only injects features to the layers of the encoder part of mainstream even corrupts the layouts from stage 1, causing performance degradation. Although ControlNet improves quality to a certain extent, this improvement is minor. In contrast, our full model realizes impressive performance gains. Note that all these settings include the proposed MAP module, thus fully demonstrating the superiority of the control branch architecture we developed. In addition, results from stage 1 are sub-optimal and even distorted, such as the hands and feet of the cartoon character. Existing generation methods struggle to rectify these flaws, but our model is capable of correcting them. For a more objective demonstration, quantitative results are given in Tab. 3, which shows our full model works better.

Ablation on MAP module.

Although the coarse 16-view videos provide layout priors, high-quality production requires more sophisticated content. To fully utilize the input monocular video, we design the MAP attention

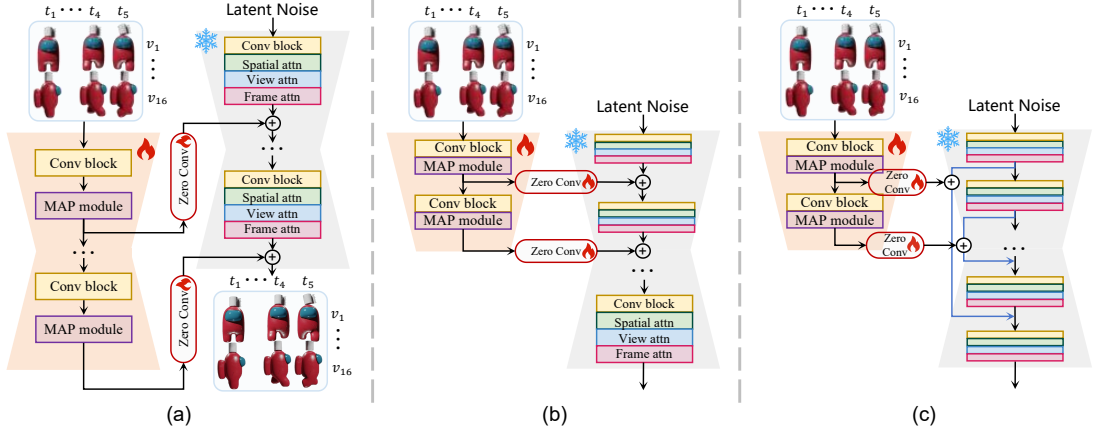


Fig. 8: Network architectures of different condition branches. We show three optional structure condition branches used in our second stage: (a) our method, (b) T2I-adapter, and (c) ControlNet.

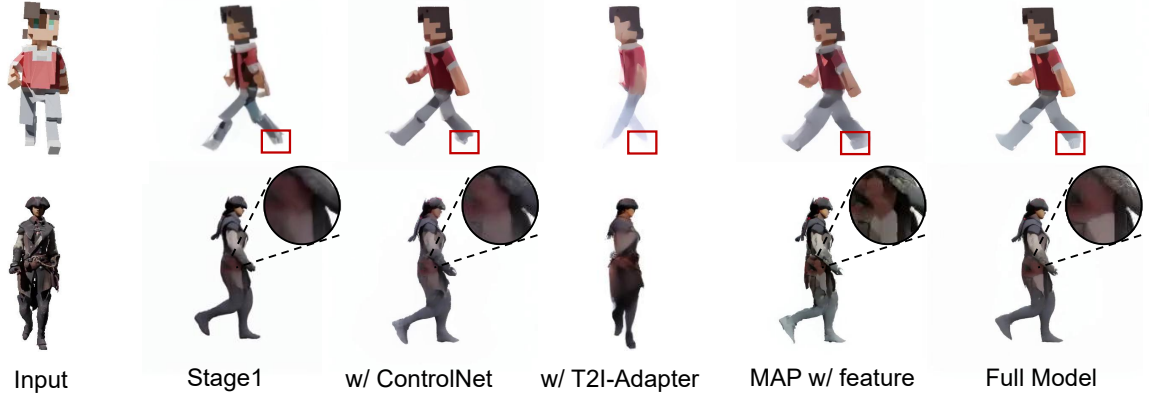


Fig. 9: Visual results of ablation study. Results from stage 1 are rough and only serve as a hint for subsequent generation. Compared with mainstream condition branches such as ControlNet or T2i-adapter, the branch we developed performs best. We also adapt the MAP module to take the model feature as the \mathbf{Q} matrix, which introduces undesired color offset and artifacts.

module to combine its high-quality feature with multi-view layouts. Considering that the coarse results contain 16 views but the reference video only records one viewpoint, we elaborate a novel cross attention to fuse them. Conventional operations [57, 58] usually choose features of the current model as \mathbf{Q} matrix, and calculate \mathbf{KV} matrices based on the external condition. We initially used a similar calculation but found that the trained model tends to introduce obvious visual artifacts as shown in “MAP w/ feature” of Fig. 9. Therefore, we take the high-quality reference video as \mathbf{Q} matrix and use coarse features to calculate the \mathbf{KV} matrices as the conditions in our

branch. Experiments demonstrate that this operation works better in our pipeline, bringing better performance. Tab. 3 shows that compared with “MAP w/ feature”, our full model achieves better results. These experiments demonstrate that compared with upscaling multi-view content based on a single view, extending high-quality reference to novel views based on the corresponding layouts is more efficient.

6 Limitations and Conclusion

Limitation

We develop 4DVD for dynamic 3D objects generation based on monocular videos, and achieve better results than previous works. However, there are still some challenges to be addressed. The current model may not perform well in complex cases with detailed geometry and complex motions due to the capacity of the network and training data. Nevertheless, 4DVD can serve as a basic model for future research, and the generation of challenging cases can be improved by incorporating more delicate designs and high-quality training data. We believe the same method can also be applied to diffusion transformers [59, 60] to address complex situations better and leave it as future work.

In conclusion, we introduce 4DVD, a cascaded latent video diffusion model designed for multi-view video generation and 4D creation, which could generate temporally consistent and high-resolution videos from a monocular video input. Previous methods attempted to directly model high-dimensional data such as 4D. In contrast, we achieve much more efficient performance by decoupling 4D generation into multi-view layout prediction and structure-aware conditional generation. Benefits from its unprecedented dense view output, the produced results remain high consistency across time and views, while explicit 4D representations, such as 4D Gaussian, can be effectively reconstructed. The model learns temporal consistent multi-view layouts by training simultaneously with 16 views, and utilizes the curated MAP module to achieve high-quality conditional spatio-temporal generation. Extensive experiments demonstrate 4DVD enables more multi-view and temporally consistent 4D generation than existing methods. As an attempt at automatic 4D creation, we believe the proposed cascaded idea provides valuable insight for future research and serves as a basic model for 4D content generation. Moreover, the proposed MAP module enables cross-view consistent conditional multi-view generation. We believe this can also be applied to scene-level generation and leave it as our future work.

Acknowledgments. This work is supported by Guangdong Provincial Key Laboratory of

Ultra High Definition Immersive Media Technology (Grant No. 2024B1212010006).

Data Availability. The Objaverse dataset used in this paper is available at <https://objaverse.allenai.org/>.

References

- [1] Jiang, Y., Zhang, L., Gao, J., Hu, W., Yao, Y.: Consistent4d: Consistent 360° dynamic object generation from monocular video. In: The Twelfth International Conference on Learning Representations (2024)
- [2] Singer, U., Sheynin, S., Polyak, A., Ashual, O., Makarov, I., Kokkinos, F., Goyal, N., Vedaldi, A., Parikh, D., Johnson, J., Taigman, Y.: Text-to-4d dynamic scene generation. In: Proceedings of the 40th International Conference on Machine Learning (2023)
- [3] Bahmani, S., Skorokhodov, I., Rong, V., Wetzstein, G., Guibas, L., Wonka, P., Tulyakov, S., Park, J.J., Tagliasacchi, A., Lindell, D.B.: 4d-fy: Text-to-4d generation using hybrid score distillation sampling. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2024)
- [4] Ling, H., Kim, S.W., Torralba, A., Fidler, S., Kreis, K.: Align your gaussians: Text-to-4d with dynamic 3d gaussians and composed diffusion models. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2024)
- [5] Ren, J., Pan, L., Tang, J., Zhang, C., Cao, A., Zeng, G., Liu, Z.: Dreamgaussian4d: Generative 4d gaussian splatting. arXiv preprint arXiv:2312.17142 (2024)
- [6] Xie, Y., Yao, C.-H., Voleti, V., Jiang, H., Jampani, V.: Sv4d: Dynamic 3d content generation with multi-frame and multi-view consistency. arXiv preprint arXiv:2407.17470 (2024)
- [7] Ren, J., Xie, K., Mirzaei, A., Liang, H., Zeng, X., Kreis, K., Liu, Z., Torralba, A., Fidler, S., Kim, S.W., Ling, H.: L4gm: Large 4d gaussian reconstruction model. arXiv preprint arXiv:2406.10324 (2024)

- [8] Yu, H., Wang, C., Zhuang, P., Menapace, W., Siarohin, A., Cao, J., Jeni, L.A., Tulyakov, S., Lee, H.-Y.: 4real: Towards photorealistic 4d scene generation via video diffusion models. In: *Advances in Neural Information Processing Systems* (2024)
- [9] Zhang, H., Chen, X., Wang, Y., Liu, X., Wang, Y., Qiao, Y.: 4diffusion: Multi-view video diffusion model for 4d generation. *arXiv preprint arXiv:2405.20674* (2024)
- [10] Deitke, M., Schwenk, D., Salvador, J., Weihs, L., Michel, O., VanderBilt, E., Schmidt, L., Ehsani, K., Kembhavi, A., Farhadi, A.: Objaverse: A universe of annotated 3d objects. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023)
- [11] Poole, B., Jain, A., Barron, J.T., Mildenhall, B.: Dreamfusion: Text-to-3d using 2d diffusion. In: *The Eleventh International Conference on Learning Representations* (2023)
- [12] Lin, C.-H., Gao, J., Tang, L., Takikawa, T., Zeng, X., Huang, X., Kreis, K., Fidler, S., Liu, M.-Y., Lin, T.-Y.: Magic3d: High-resolution text-to-3d content creation. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023)
- [13] Wang, Z., Lu, C., Wang, Y., Bao, F., Li, C., Su, H., Zhu, J.: Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation. In: *Advances in Neural Information Processing Systems* (2023)
- [14] Yu, W., Yuan, L., Cao, Y.-P., Gao, X., Li, X., Hu, W., Quan, L., Shan, Y., Tian, Y.: Hifi-123: Towards high-fidelity one image to 3d content generation. *arXiv preprint arXiv:2310.06744* (2024)
- [15] Weng, H., Yang, T., Wang, J., Li, Y., Zhang, T., Chen, C.L.P., Zhang, L.: Consistent123: Improve consistency for one image to 3d object synthesis. *arXiv preprint arXiv:2310.08092* (2023)
- [16] Yi, T., Fang, J., Wang, J., Wu, G., Xie, L., Zhang, X., Liu, W., Tian, Q., Wang, X.: Gaussiandreamer: Fast generation from text to 3d gaussians by bridging 2d and 3d diffusion models. *arXiv preprint arXiv:2310.08529* (2024)
- [17] Chen, Z., Wang, F., Wang, Y., Liu, H.: Text-to-3d using gaussian splatting. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2024)
- [18] Sun, J., Zhang, B., Shao, R., Wang, L., Liu, W., Xie, Z., Liu, Y.: Dreamcraft3d: Hierarchical 3d generation with bootstrapped diffusion prior. In: *The Twelfth International Conference on Learning Representations* (2024)
- [19] Tang, J., Ren, J., Zhou, H., Liu, Z., Zeng, G.: Dreamgaussian: Generative gaussian splatting for efficient 3d content creation. In: *The Twelfth International Conference on Learning Representations* (2024)
- [20] Liang, Y., Yang, X., Lin, J., Li, H., Xu, X., Chen, Y.: Luciddreamer: Towards high-fidelity text-to-3d generation via interval score matching. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2024)
- [21] Zhou, L., Shih, A., Meng, C., Ermon, S.: Dreampropeller: Supercharge text-to-3d generation with parallel sampling. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2024)
- [22] Yang, S., Wang, Y., Li, H., Meng, J., Wu, Y., Meng, X., Zhang, J.: Hybrid fourier score distillation for efficient one image to 3d object generation. *arXiv preprint arXiv:2405.20669* (2024)
- [23] Hong, Y., Zhang, K., Gu, J., Bi, S., Zhou, Y., Liu, D., Liu, F., Sunkavalli, K., Bui, T., Tan, H.: LRM: Large reconstruction model for single image to 3d. In: *The Twelfth International Conference on Learning Representations* (2024)
- [24] Jiang, H., Jiang, Z., Zhao, Y., Huang, Q.: LEAP: Liberate sparse-view 3d modeling

- from camera poses. In: The Twelfth International Conference on Learning Representations (2024)
- [25] Wang, P., Tan, H., Bi, S., Xu, Y., Luan, F., Sunkavalli, K., Wang, W., Xu, Z., Zhang, K.: PF-LRM: Pose-free large reconstruction model for joint pose and shape prediction. In: The Twelfth International Conference on Learning Representations (2024)
- [26] Zou, Z.-X., Yu, Z., Guo, Y.-C., Li, Y., Liang, D., Cao, Y.-P., Zhang, S.-H.: Triplane meets gaussian splatting: Fast and generalizable single-view 3d reconstruction with transformers. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2024)
- [27] Wei, X., Zhang, K., Bi, S., Tan, H., Luan, F., Deschaintre, V., Sunkavalli, K., Su, H., Xu, Z.: Meshlrn: Large reconstruction model for high-quality mesh. arXiv preprint arXiv:2404.12385 (2024)
- [28] Tochilkin, D., Pankratz, D., Liu, Z., Huang, Z., Letts, A., Li, Y., Liang, D., Laforte, C., Jampani, V., Cao, Y.-P.: Tripopr: Fast 3d object reconstruction from a single image. arXiv preprint arXiv:2403.02151 (2024)
- [29] Tang, J., Chen, Z., Chen, X., Wang, T., Zeng, G., Liu, Z.: Lgm: Large multi-view gaussian model for high-resolution 3d content creation. In: European Conference on Computer Vision (2025)
- [30] Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: representing scenes as neural radiance fields for view synthesis. *Commun. ACM* (2021)
- [31] Kerbl, B., Kopanas, G., Leimkühler, T., Dretakis, G.: 3d gaussian splatting for real-time radiance field rendering. arXiv preprint arXiv:2308.04079 (2023)
- [32] Shi, Y., Wang, P., Ye, J., Mai, L., Li, K., Yang, X.: MVDream: Multi-view diffusion for 3d generation. In: The Twelfth International Conference on Learning Representations (2024)
- [33] Wang, P., Shi, Y.: Imagedream: Image-prompt multi-view diffusion for 3d generation. arXiv preprint arXiv:2312.02201 (2023)
- [34] Long, X., Guo, Y.-C., Lin, C., Liu, Y., Dou, Z., Liu, L., Ma, Y., Zhang, S.-H., Habermann, M., Theobalt, C., Wang, W.: Wonder3d: Single image to 3d using cross-domain diffusion. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2024)
- [35] Li, P., Liu, Y., Long, X., Zhang, F., Lin, C., Li, M., Qi, X., Zhang, S., Luo, W., Tan, P., Wang, W., Liu, Q., Guo, Y.: Era3d: High-resolution multiview diffusion using efficient row-wise attention. arXiv preprint arXiv:2405.11616 (2024)
- [36] Zhao, Y., Yan, Z., Xie, E., Hong, L., Li, Z., Lee, G.H.: Animate124: Animating one image to 4d dynamic scene. arXiv preprint arXiv:2311.14603 (2024)
- [37] Yin, Y., Xu, D., Wang, Z., Zhao, Y., Wei, Y.: 4dgen: Grounded 4d content generation with spatial-temporal consistency. arXiv preprint arXiv:2312.17225 (2024)
- [38] Zeng, Y., Jiang, Y., Zhu, S., Lu, Y., Lin, Y., Zhu, H., Hu, W., Cao, X., Yao, Y.: Stag4d: Spatial-temporal anchored generative 4d gaussians. In: European Conference on Computer Vision (2025)
- [39] Pan, Z., Yang, Z., Zhu, X., Zhang, L.: Efficient4d: Fast dynamic 3d object generation from a single-view video. arXiv preprint arXiv:2401.08742 (2024)
- [40] Yang, Z., Pan, Z., Gu, C., Zhang, L.: Diffusion²: Dynamic 3d content generation via score composition of video and multi-view diffusion models. arXiv preprint arXiv:2404.02148 (2024)
- [41] Pumarola, A., Corona, E., Pons-Moll, G., Moreno-Noguer, F.: D-nerf: Neural radiance fields for dynamic scenes. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2021)

- [42] Cao, A., Johnson, J.: Hexplane: A fast representation for dynamic scenes. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2023)
- [43] Wu, G., Yi, T., Fang, J., Xie, L., Zhang, X., Wei, W., Liu, W., Tian, Q., Wang, X.: 4d gaussian splatting for real-time dynamic scene rendering. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2024)
- [44] Liang, H., Yin, Y., Xu, D., Liang, H., Wang, Z., Plataniotis, K.N., Zhao, Y., Wei, Y.: Diffusion4d: Fast spatial-temporal consistent 4d generation via video diffusion models. arXiv preprint arXiv:2405.16645 (2024)
- [45] Sun, W., Chen, S., Liu, F., Chen, Z., Duan, Y., Zhang, J., Wang, Y.: Dimensionx: Create any 3d and 4d scenes from a single image with controllable video diffusion. arXiv preprint arXiv:2411.04928 (2024)
- [46] YU, M., Hu, W., Xing, J., Shan, Y.: Trajectorycrafter: Redirecting camera trajectory for monocular videos via diffusion models. arXiv preprint arXiv:2503.05638 (2025)
- [47] Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2022)
- [48] Wan, T., Wang, A., Ai, B., Wen, B., Mao, C., Xie, C.-W., Chen, D., Yu, F., Zhao, H., Yang, J., Zeng, J., Wang, J., Zhang, J., Zhou, J., Wang, J., Chen, J., Zhu, K., Zhao, K., Yan, K., Huang, L., Feng, M., Zhang, N., Li, P., Wu, P., Chu, R., Feng, R., Zhang, S., Sun, S., Fang, T., Wang, T., Gui, T., Weng, T., Shen, T., Lin, W., Wang, W., Wang, W., Zhou, W., Wang, W., Shen, W., Yu, W., Shi, X., Huang, X., Xu, X., Kou, Y., Lv, Y., Li, Y., Liu, Y., Wang, Y., Zhang, Y., Huang, Y., Li, Y., Wu, Y., Liu, Y., Pan, Y., Zheng, Y., Hong, Y., Shi, Y., Feng, Y., Jiang, Z., Han, Z., Wu, Z.-F., Liu, Z.: Wan: Open and advanced large-scale video generative models. arXiv preprint arXiv:2503.20314 (2025)
- [49] Yao, C.-H., Xie, Y., Voleti, V., Jiang, H., Jampani, V.: Sv4d 2.0: Enhancing spatio-temporal consistency in multi-view video diffusion for high-quality 4d generation. arXiv preprint arXiv:2503.16396 (2025)
- [50] Hu, E.J., yelong shen, Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., Chen, W.: LoRA: Low-rank adaptation of large language models. In: International Conference on Learning Representations (2022)
- [51] Voleti, V., Yao, C.-H., Boss, M., Letts, A., Pankratz, D., Tochilkin, D., Laforte, C., Rombach, R., Jampani, V.: Sv3d: Novel multi-view synthesis and 3d generation from a single image using latent video diffusion. In: European Conference on Computer Vision (2025)
- [52] Zhang, L., Rao, A., Agrawala, M.: Adding conditional control to text-to-image diffusion models. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (2023)
- [53] Mou, C., Wang, X., Xie, L., Wu, Y., Zhang, J., Qi, Z., Shan, Y.: T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models. In: Proceedings of the AAAI Conference on Artificial Intelligence (2024)
- [54] Hu, L.: Animate anyone: Consistent and controllable image-to-video synthesis for character animation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2024)
- [55] Wang, W., Lv, Q., Yu, W., Hong, W., Qi, J., Wang, Y., Ji, J., Yang, Z., Zhao, L., Song, X., Xu, J., Chen, K., Xu, B., Li, J., Dong, Y., Ding, M., Tang, J.: Cogvlm: Visual expert for pretrained language models. In: Advances in Neural Information Processing Systems (2024)
- [56] Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable effectiveness of deep features as a perceptual metric. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition

(2018)

- [57] Cao, M., Wang, X., Qi, Z., Shan, Y., Qie, X., Zheng, Y.: Masactrl: Tuning-free mutual self-attention control for consistent image synthesis and editing. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (2023)
- [58] QI, C., Cun, X., Zhang, Y., Lei, C., Wang, X., Shan, Y., Chen, Q.: Fatezero: Fusing attentions for zero-shot text-based video editing. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (2023)
- [59] Hong, W., Ding, M., Zheng, W., Liu, X., Tang, J.: Cogvideo: Large-scale pretraining for text-to-video generation via transformers. arXiv preprint arXiv:2205.15868 (2022)
- [60] Yang, Z., Teng, J., Zheng, W., Ding, M., Huang, S., Xu, J., Yang, Y., Hong, W., Zhang, X., Feng, G., Yin, D., Zhang, Y., Wang, W., Cheng, Y., Xu, B., Gu, X., Dong, Y., Tang, J.: Cogvideox: Text-to-video diffusion models with an expert transformer. arXiv preprint arXiv:2408.06072 (2025)