# A Unified Approach for Text- and Image-guided 4D Scene Generation

Yufeng Zheng[1,2,3], Xueting Li[1], Koki Nagano[1], Sifei Liu[1], Karsten Kreis[1], Otmar Hilliges[2], Shalini De Mello[1]
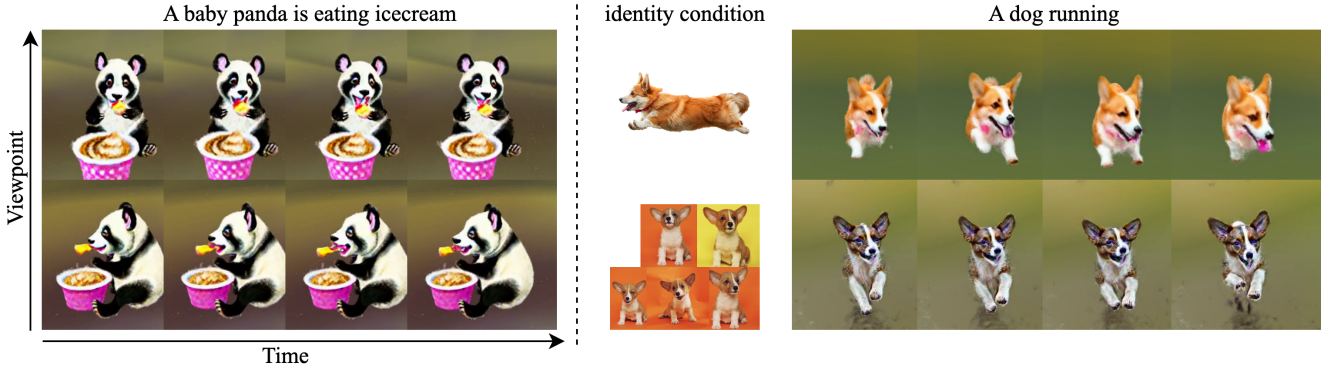[1]NVIDIA, [2]ETH Zurich, [3]Max Planck Institute for Intelligent Systems

Figure 1. **Text-to-4D.** Our method provides a unified approach for generating 4D dynamic content from a text prompt with diffusion guidance, supporting both unconstrained generation and controllable generation, where appearance is defined by one or multiple images.

## Abstract

*Large-scale diffusion generative models are greatly simplifying image, video and 3D asset creation from user-provided text prompts and images. However, the challenging problem of text-to-4D dynamic 3D scene generation with diffusion guidance remains largely unexplored. We propose Dream-in-4D, which features a novel two-stage approach for text-to-4D synthesis, leveraging (1) 3D and 2D diffusion guidance to effectively learn a high-quality static 3D asset in the first stage; (2) a deformable neural radiance field that explicitly disentangles the learned static asset from its deformation, preserving quality during motion learning; and (3) a multi-resolution feature grid for the deformation field with a displacement total variation loss to effectively learn motion with video diffusion guidance in the second stage. Through a user preference study, we demonstrate that our approach significantly advances image and motion quality, 3D consistency and text fidelity for text-to-4D generation compared to baseline approaches. Thanks to its motion-disentangled representation, Dream-in-4D can also be easily adapted for controllable generation where appearance is defined by one or multiple images, without the need to modify the motion learning stage. Thus, our method offers, for the first time, a unified approach for text-to-4D, image-to-4D and personalized 4D generation tasks.*

## 1. Introduction

The advent of large-scale text-conditioned diffusion-based generative models for images has ushered in a new era of imaginative, high-quality image synthesis [1, 29, 31]. Their simple and intuitive conditioning in the form of text prompts are game-changers in democratizing visual content creation for non-expert users. Subsequently, these developments have also led to impressive progress in (1) text- or image-conditioned static 3D content creation [15, 21, 25, 27], achieved by leveraging guidance from generic and 3D-aware [16, 17, 32] image diffusion models, and (2) video content creation via video diffusion models [2–4, 33].

However, for various real-world applications such as gaming, AR/VR, and advertising, synthesizing static 3D assets alone does not suffice. It is desirable to also animate 3D assets using intuitive user-provided text prompts to further save animators' time and level-of-expertise. To go beyond static 3D content creation, we delve into the largely unexplored problem of text-conditioned 4D scene generation, a.k.a., *text-to-4D* synthesis with diffusion guidance. This is a challenging problem encompassing both text-to-3D and text-to-video synthesis. It requires learning not only a 3D-consistent representation of a static scene capable of free-view rendering, but also its plausible and semantically-correct dynamic 3D motion over time.

---

Work was done during an internship at NVIDIA.

(a) Multi-view images of a generated 3D asset without 3D diffusion guidance. It suffers from the Janus problem.

(b) The presence of the Janus problem, further hinders learning of the cape's correct temporal motion (top row).

Figure 2. We show the importance of a high-quality static asset for 4D content generation. The prompt for this scene is 'Superhero dog with red cape flying through the sky'.

The concurrent pioneering work [34], constitutes the sole existing attempt to address this problem. It proposes a two-stage approach: the first to learn a static 3D asset, and the second to optimize its full dynamic representation with guidance from a video diffusion model [33]. It further models the dynamic representation via a neural hexplane [5]. While impressive in demonstrating feasibility, this early work leaves much room for improvement in terms of robustness, quality and realism. Additionally, it does not solve the problems of image-to-4D or personalized-4D content creation, wherein, in addition to a text prompt, an image or a set of images is provided as input to control the appearance of the 4D outcome.

To address these challenges, we propose a novel method for text-to-4D dynamic scene synthesis, named *Dream-in-4D*. It employs a two-stage approach, to first learn a static scene representation and then its motion. Our first primary insight is that achieving high-quality, 3D-consistent static reconstruction in the first stage is crucial for successfully learning motion in the second stage. For example, in Fig. 2 we show a a multi-view *inconsistent* 3D dog with two heads due to the Janus problem learned in the first stage. It introduces significant ambiguity for the second dynamic stage and substantially undermines the quality of the learned motion (in the dog's cape). However, relying solely on guidance from image or video diffusion models for static text-to-3D synthesis, as proposed in [34], easily encounters the Janus problem (see the first row of Fig. 4). Therefore, we leverage 3D-aware [16, 32] and standard image diffusion models [1, 29] in stage-one to achieve high-quality view-consistent text-to-3D synthesis, along with video diffusion guidance [3] in stage-two to learn realistic motion. This forms our first key contribution, which is to leverage guidance from a carefully-designed combination of image, 3D, and video diffusion models to effectively solve the task of text-to-4D synthesis.

Our approach is further motivated by the observation that fine-tuning the static model with video diffusion models in stage-two leads to lower visual quality and prompt fidelity, primarily because these models are trained with lower quality videos compared to image models. To address this problem, our insight is to decompose the synthesis process into two *distinct* training stages, the first of which is designed to learn a high-quality static 3D asset and the second dedicated to effectively animating it with the provided text prompt, while keeping the pre-trained static 3D asset unchanged. This, in turn, requires a 4D neural representation that fully disentangles the canonical static representation and its motion. However, this cannot be achieved with the hexplane [5] representation proposed in [34], which entangles the static representation and its motion. To this end, we propose to use a variant of a deformable neural radiance field (D-NeRF) [26] for the task of 4D content generation. A D-NeRF consists of a canonical 3D NeRF [22] and a 4D deformation MLP that maps time-dependent deformed space to the common canonical static space. With the proposed disentangled representation for 4D scene synthesis, we can freeze the pre-trained high-quality static 3D model from stage-one and only optimize the deformation field using video diffusion guidance in stage-two. To successfully learn detailed and realistic motion, we further encode the 4D deformation field with multi-resolution feature grids and regularize motion using a novel total variation loss on the rendered displacement maps. We find that the former enhances detailed motion while the latter reduces spatial and temporal jitter.

Through a user preference study on diverse text prompts, we show that our algorithm achieves significant improvements in visual quality, 3D consistency, prompt matching and motion quality compared to alternate baselines. Furthermore, the ability to disentangle the canonical and motion representation allows for easy adaptation to image-conditioned 4D generation, without requiring modifications to the motion learning stage. Thus, we demonstrate, for the first time, image-to-4D generation given a single-view image, and personalized 4D generation using 4-6 casually captured images of a subject (See Fig. 1) with our unified *Dream-in-4D* method.

In summary, our key contributions include:

1. We propose to combine image, 3D-aware and video diffusion priors for the text-to-4D task, significantly improving the visual quality, 3D consistency and text-fidelity of the learned static assets in the first stage.
2. By explicitly disentangling the static representation from its deformation, our method preserves the high quality static asset during motion learning.
3. We propose to use a multi-resolution feature grid and a total variation loss on the deformation field to effectively learn motion with video diffusion guidance.
4. We demonstrate that our method offers, for the first time, a unified approach for text-to-4D, image-to-4D and personalized 4D generation tasks.

## 2. Related Work

**Dynamic neural radiance field.** Modeling dynamic 3D content with NeRFs has been extensively studied in the novel-view synthesis literature. To extend NeRFs to dynamic scene modeling, previous works either learn a high-dimensional radiance field conditioned on temporal embeddings [14, 19], or a separate deformation mapping to model motion [24, 26]. To speed up training and inference, plane- and voxel-based feature grids are combined with MLPs to formulate efficient hybrid NeRF representations [6, 11, 23], which are extended to dynamic scene modeling by learning additional planes for the temporal dimension [5, 12]. In this work, we leverage a deformable NeRF representation where the canonical geometry and 4D deformation field are both encoded by multi-resolution feature grids [23]. As a result, the geometry and motion are fully disentangled, which not only eases motion learning but also allows easy adaptation to various applications such as image(s)-to-4D video generation.

**Diffusion models.** Recently, diffusion models have revolutionized the computer vision community by showing remarkable advancements in image, video or novel view image synthesis. Seminal works such as Stable Diffusion [29] and DeepFloyd [1] take a text prompt as input and produce high-quality images that align with the prompt. Leveraging large-scale image datasets, these diffusion models learn various prior knowledge ranging from object appearance to complex scene layout. One line of subsequent works [2–4] fine-tune text-to-image diffusion models on video datasets, successfully extending them to generate realistic videos matching both the object and motion described by the input prompt. Another line of methods [16, 17] learn 3D-aware diffusion models with images rendered from synthetic objects [8, 9]. By conditioning on camera parameters, these methods produce novel view images of an object that are consistent with each other and align with the observed view.

**Text/image(s) to 3D with diffusion priors.** Beyond direct sampling from the diffusion models, several works employ image diffusion priors as an optimization signal for 3D generation. The pioneering work, DreamFusion [25], optimizes a 3D model by presenting its renderings to a text-to-image diffusion model and acquiring gradient supervision through Score Distillation Sampling (SDS). Subsequent works enhance the synthesis quality and speed by incorporating the mesh representation [15, 36], advancing score distillation [37–39], exploring representations in the latent space [21], or disentangling geometry and texture [7]. Yet, these methods suffer from the Janus problem due to the lack of 3D prior in the text-to-image diffusion models. To overcome this limitation, several works [13, 16, 17, 32] leverage 3D-aware diffusion models as supervisions, generating 3D objects that are consistent across multiple views. In addition to text, a few

approaches further take one or multiple images as inputs and reconstruct a 3D object matching both the prompt and the image(s). The former [20, 27, 35] simultaneously utilize text-to-image and novel view diffusion models [16], while the latter [28, 32] overfit a diffusion model on a few images depicting the same subject to achieve personalized diffusion guidance.

**Text to 4D with diffusion priors.** In this paper, we go beyond 2D/3D generation and aim to synthesize a 4D video given a text prompt. This is hitherto a highly challenging and under-explored domain. The most relevant work to ours is MAV3D [34], which optimizes a 4D scene by leveraging a pre-trained video diffusion model. However, due to the entanglement of geometry and motion, as well as a lack of 3D-aware prior, this method suffers from the Janus problem and produces low quality texture. To resolve these limitations, we leverage a carefully designed combination of image, video and 3D-aware diffusion models and fully disentangle the geometry and motion. Our method synthesizes multi-view consistent 4D videos with realistic appearance and motion. Furthermore, the disentanglement of canonical and motion representations readily enables novel applications such as image-to-4D and personalized 4D generation.

## 3. Method

Given a text prompt and optionally one or a few images to specify the object appearance, we aim to generate a 4D video that matches both the object and the motion described in the prompt. To this end, we propose a two-stage training pipeline. In the first static stage (Sec. 3.1), we synthesize a high-quality static 3D scene using both 2D and 3D diffusion priors. In the second dynamic stage (Sec. 3.2), we learn the 3D motion of the scene using a video diffusion model, while keeping the static scene representation intact.

### 3.1. Static Stage

The goal of the static stage is to generate a high-quality 3D scene that aligns with the text prompt. To efficiently learn a canonical model of a 3D scene, we opt for the NeRF [22] representation with multi-resolution hash-encoded features [23], which is extensively used in previous text-to-3D methods [15, 21, 25, 38, 39]. This static 3D model is accompanied by a deformation field to represent the dynamic motion of the 3D scene in the subsequent motion learning stage. Two elements of the static stage are crucial for 3D asset quality and play an important role in facilitating motion learning in the subsequent dynamic stage: (1) the generated 3D object(s) should be view-consistent (i.e., free of the Janus problem) and (2) should follow the spatial composition described in the text prompt. Intuitively, the former reduces contradictory gradients from different views in deformation optimization while the lat-
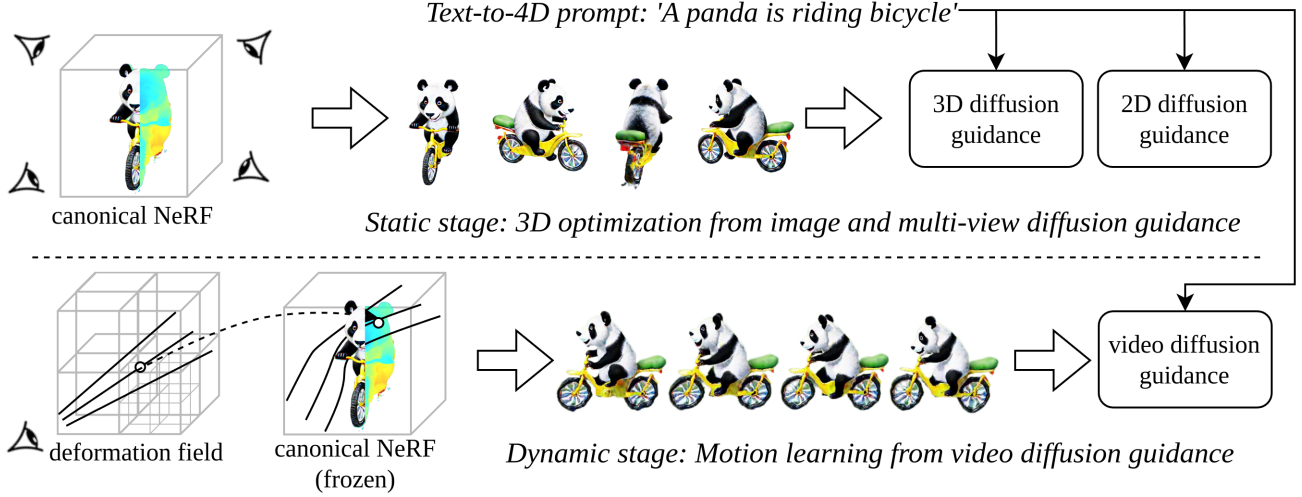
Figure 3. **Method overview.** Adopting a two-stage approach, *Dream-in-4D* first utilizes 3D and 2D diffusion guidance to learn a static 3D asset based on the provided text prompt (top). Then, it optimizes a deformation field using video diffusion guidance to model the motion described in the text prompt (bottom). Featuring a motion-disentangled D-NeRF representation, our method freezes the pre-trained static canonical asset while optimizing for the motion, achieving high quality view-consistent 4D dynamic content with realistic motion.



Figure 4. **Static stage.** Without StableDiffusion guidance, the learned static model fails to learn the correct composition. Without MVDream guidance, the learned assets suffer from the Janus problem and contain multiple faces. Using guidance from both StableDiffusion and MVDream, results in the best text prompt fidelity and 3D consistency.

ter eases motion learning by presenting a reasonable spatial layout of multiple objects. For instance, the panda sitting on top of a bike in Fig. 3 sets a good starting point for the dynamic stage to learn the "riding" motion. To achieve these goals, we propose to utilize both 3D and generic 2D diffusion models for the static stage. This is in spirit similar to prior work [27] for image-to-3D synthesis. In the following, we introduce the 3D and 2D diffusion guidance used for stage-one.

3D diffusion models [16, 17, 32] take camera parameters with a text prompt or an image as inputs and synthesize novel view images of the target object. Fine-tuned from image diffusion models with rendered images of synthetic 3D data [10], 3D diffusion models provide valuable prior knowledge of the 3D world and enforce different views of a 3D object to be consistent. For text to 3D generation, we adopt the MVDream [32] model to provide a 3D prior. Specifically, we render the synthesized 3D object from four different viewpoints (i.e., front, back, and two side views) and obtain guidance from a pre-trained MVDream model through the SDS loss [25]. We use a reconstruction formulation of the SDS loss, similar to [32]. The 3D guidance loss is denoted as $\mathcal{L}_{3D}(I)$.

However, due to the limited scale and synthetic nature of the 3D training datasets, static Nerf models optimized using MVDream alone tend to have synthetic-looking texture [32], and occasionally fail to produce realistic scene layouts (see the second row of Fig. 4 where objects in the prompt are missing from the scene). Meanwhile, we observe that image diffusion models trained with large-scale 2D images encourage both realistic appearance and reasonable scene layouts, but by themselves, easily suffer from the Janus problem (see the first row of Fig. 4). Thus, we propose to combine 2D diffusion guidance with 3D guidance in stage-one. Specifically, we use StableDiffusion-v2.1 with the SDS objective (denoted as $\mathcal{L}_{2D}(I)$) to provide 2D guidance. The overall objective for stage-one can be expressed as:

$$\mathcal{L} = \lambda_{2D}\mathcal{L}_{2D}(I) + \lambda_{3D}\mathcal{L}_{3D}(I),$$

where $I$ denotes the set of rendered images from the sampled camera viewpoints, and $\lambda_{2D/3D}$ are the weights for the 2D and 3D guidances (see Supplement for the loss weights).

4

Hexplane static      Hexplane dynamic      Ours dynamic

Figure 5. **Hexplane v.s. deformable NeRF.** With a hexplane representation, even though the static stage successfully learns a high-quality 3D asset (column 1), the motion learning stage with video diffusion guidance still leads to degradation in texture and re-appearance of the Janus problem (column 2).

As shown in Fig. 4, by combining both the 2D and 3D diffusion models for guidance, our static stage generates 3D-consistent object(s) with realistic texture and plausible scene layouts.

## 3.2. Dynamic Stage

In the dynamic stage, our goal is to learn a deformation field that animates the 3D scene generated in the static stage using guidance from a video diffusion model. As aforementioned, our key observation is that although video diffusion models provide a valuable motion prior, they are not 3D-aware and tend to produce unappealing visual results (see Fig. 5, column 2). Therefore, we propose to fully disentangle the static model and the motion by freezing the NeRF network learned in the static stage and only learn the deformation field to match the motion described in the text prompt in the dynamic stage. Such a design brings two advantages: (1) it preserves the view consistency and high-quality texture learned in the static stage and (2) it readily enables applications such as image-to-4D and personalized 4D generation (see Sec. 3.3).

**Motion-disentangled 4D representation.** Our dynamic 4D representation consists of a canonical 3D radiance field (as described in Sec. 3.1) and a deformation field. The deformation field is a 4D to 3D time-dependent mapping $D(\mathbf{x_d}, t) \rightarrow \mathbf{x_c}$, where $\mathbf{x_d}$ is a 3D point's location in deformed space at time $t$, and $\mathbf{x_c}$ is its corresponding canonical location. Our insight is that the deformable field should be smooth both spatially and temporally due to the limited elasticity and velocity of the object. As a result, the deformation field does not require as high-resolution a feature grid as its static canonical 3D counterpart. Therefore, we utilize a 4D multi-resolution hash-encoded feature grid with a maximum resolution of 232 for the deformation field, in contrast to the maximum resolution of 4096 for the canonical static NeRF representation. Additionally, we found the usage of multi-resolution features to be crucial for learning correct local motion (see Fig. 6).

**Motion optimization with video diffusion models.** The deformation field is optimized via score distillation sampling using a video diffusion model. Specifically,

we sample a static camera parameter, and render a 24-frame video $\mathbf{V}$ from our 4D representation. The time stamps are sampled evenly and the length of the video is randomly chosen between 0.8 and 1 (assuming that the full length is 1). We leverage a variant of the SDS loss [39] for the video diffusion guidance, where we predict the original video with 1-step denoising, and use a combination of a latent feature loss and a decoded RGB space loss. The video diffusion guidance loss can be expressed as $\mathcal{L}_{video}(\mathbf{V}) = \mathcal{L}_{latent}(\mathbf{V}) + \lambda_{dec}\mathcal{L}_{dec}(\mathbf{V})$, where $\lambda_{dec} = 0.1$. We choose to use the Zeroscope [3] video diffusion model as our motion prior, but our method is robust to other models such as Modelscope [2] (see our supplementary video for results). We found that matching the training resolution of video diffusion models to be important for successfully distilling motion priors. For Zeroscope, we render videos at a resolution of $144 \times 72$ and upsample them to $576 \times 320$ when training with video diffusion guidance.

**Total variation motion regularization.** To reduce temporal and spatial jitter in motion, we propose to use a novel total variation loss for the learned deformation (see Fig. 6). Specifically, in addition to the RGB video $\mathbf{V}$, we also render a video for the 3D displacements $\mathbf{D}$. The total variation loss on the rendered displacement video $\mathbf{D}$ can be expressed as:

$$\mathcal{L}_{TV}(D) = \sum_{x,y,t} (||\mathbf{D}_{x-1,y,t} - \mathbf{D}_{x,y,t}||_2^2$$
$$+ ||\mathbf{D}_{x,y-1,t} - \mathbf{D}_{x,y,t}||_2^2 + ||\mathbf{D}_{x,y,t-1} - \mathbf{D}_{x,y,t}||_2^2).$$

The overall objective function for the second stage is then:

$$\mathcal{L} = \mathcal{L}_{video}(\mathbf{V}) + \lambda_{TV}\mathcal{L}_{TV}(\mathbf{D}),$$

where $\lambda_{TV} = 1000$.

## 3.3. 4D Generation Given One or Multiple Images

While text-to-4D generation is useful for many scenarios, there is a common desire to create content that features a specific object. However, language alone may be insufficient to describe the unique appearance of a given object. Thanks to the full disentanglement of its static and dynamic parts, our method can be easily extended to image-guided 4D generation, without modifying the motion learning stage. In the following, we show that this can be done by simply replacing the diffusion models used in the static stage of our method.

**Image-to-4D generation.** Given a single image, we reconstruct the corresponding 3D asset by replacing MVDream with an image-conditioned 3D diffusion model. Specifically, we use zero123-xl [16] as our 3D diffusion model and DeepfloydIF [1] as our 2D diffusion model. Additionally, we supervise the reference view with the
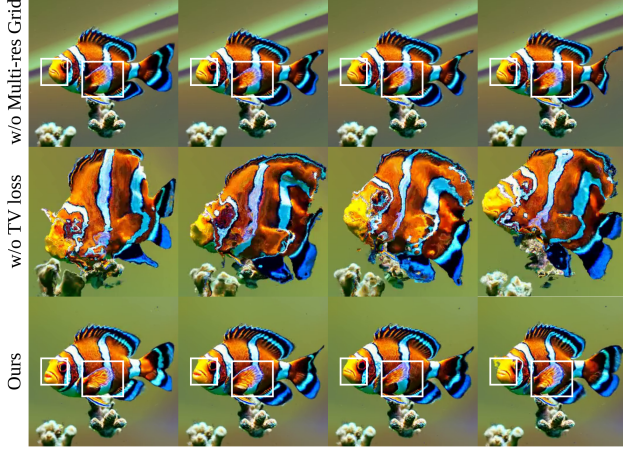
5

Figure 6. **Deformation learning.** The deformation MLP equipped with positional encoding instead of multi-resolution feature grids cannot capture local motions (mouth and fin in row 1). Without the proposed total variation loss on the displacement, the learned deformation contains substantial noise (row 2). Our approach, with both, results in the best quality.

given image and its estimated foreground mask, similar to [16, 27]. Fig. 8 shows examples of synthesized 4D videos from a single-view image.

**Personalized 4D generation.** Given a few casually captured images of an object, Dreambooth [30] finetunes image diffusion models to generate personalized images of this object given a text prompt. By replacing our generic image diffusion model with a finetuned, personalized version, we can create personalized 3D assets given a text prompt and a few casual images. Specifically, we use personalized StableDiffusion together with MVDream for this task. We show synthesized 4D videos in Fig. 9.

## 4. Results

### 4.1. Text-to-4D Generation

In Fig. 7, we show qualitative results of our method on text-to-4D generation for various text prompts. Since there is no publicly available implementation of prior work [34], we compare our method against several ablative baselines, qualitatively and with a user preference study. Video results of all methods are provided in the supplementary material for better assessment of the motion.

**User study.** We carry out a user preference study to evaluate quality along the dimensions of (1) alignment to the input text prompt, (2) motion quality and (3) 3D consistency and visual quality. In total, we gather 540 votes from 18 users. For each vote, we present the participant with results from our method as well as from other ablative baselines, and ask the participant to pick the best method given one of the three evaluation metrics described above. The results are shown

in Tab. 1.

| Metric | w/o 2D guidance | Ours |
| --- | --- | --- |
| text alignment | 11.67% | 88.33% |

(a) **Text alignment ablation.** Without 2D diffusion guidance, the learned 3D asset might fail to generate all the required components of a scene or fail to produce a plausible layout (see Fig. 4).

| Metric | w/o Multi-res Grid | w/o TV loss | Ours |
| --- | --- | --- | --- |
| motion quality | 42.22% | 2.78% | 55.00% |

(b) **Motion quality ablation.** Without the multi-resolution feature grid, detailed local motion cannot be learned. Without the proposed TV loss, the generated motion contains substantial noise (see Fig. 6).

| Metric | w/o 3D guidance | Hexplane | Ours |
| --- | --- | --- | --- |
| visual&3D | 6.12% | 0.00% | 93.88% |

(c) **Visual quality and 3D consistency abation.** Without 3D diffusion guidance in the first stage, the learned 3D assets suffer from the Janus problem (see Fig. 4). With a hexplane 4D representation, the learned high-quality 3D asset from the static stage cannot be preserved in the dynamic stage, leading to lower visual quality and re-appearance of the Janus problem (see Fig. 5).

Table 1. **User study results.** Our method is preferred over all ablative baselines in terms of (a) text alignment, (b) motion quality and (c) visual quality and 3D consistency.

**3D and 2D diffusion guidance in the static stage.** Our method combines 3D and 2D diffusion guidance to learn the static model. Fig. 4 shows qualitative comparisons, ablating the guidance used in the static stage. Without 2D diffusion guidance, the method often fails to produce the correct layout of the scene, and sometimes produces synthetic-looking texture. This is also reflected by the user study in Tab. 1a. Without 3D diffusion guidance, the learned assets suffer from severe Janus problems and do not have plausible shapes (see row 1 of Fig. 4 and 'w/o 3D guidance' in Tab. 1c). By combining both 3D and 2D guidance, our method reconstructs 3D-consistent static scenes with plausible compositions and realistic textures.

**Deformation field and motion regularization.** To learn better motion, we propose to use a multi-resolution hash-encoded 4D feature grid for the deformation MLP. We ablate this choice against a baseline MLP with positional encoding [22]. In Fig. 6, our method learns more local motion around the mouth and fin area of the clown fish. We also ablate the total variation (TV) loss on the displacement map and show that the learned motion presents substantial noise when not using the TV loss (Fig. 6). These observations about the learned motion quality are also reflected by the user study in Tab. 1b.

**D-NeRF v.s. hexplane representation** We also ablate our deformable NeRF representation against the hexplane 4D

A monkey is eating candy bar

A man is drinking beer

A baby is eating ice cream

A fox is playing video game

A cat is singing

A goat is drinking beer

Emoji of a baby panda reading book

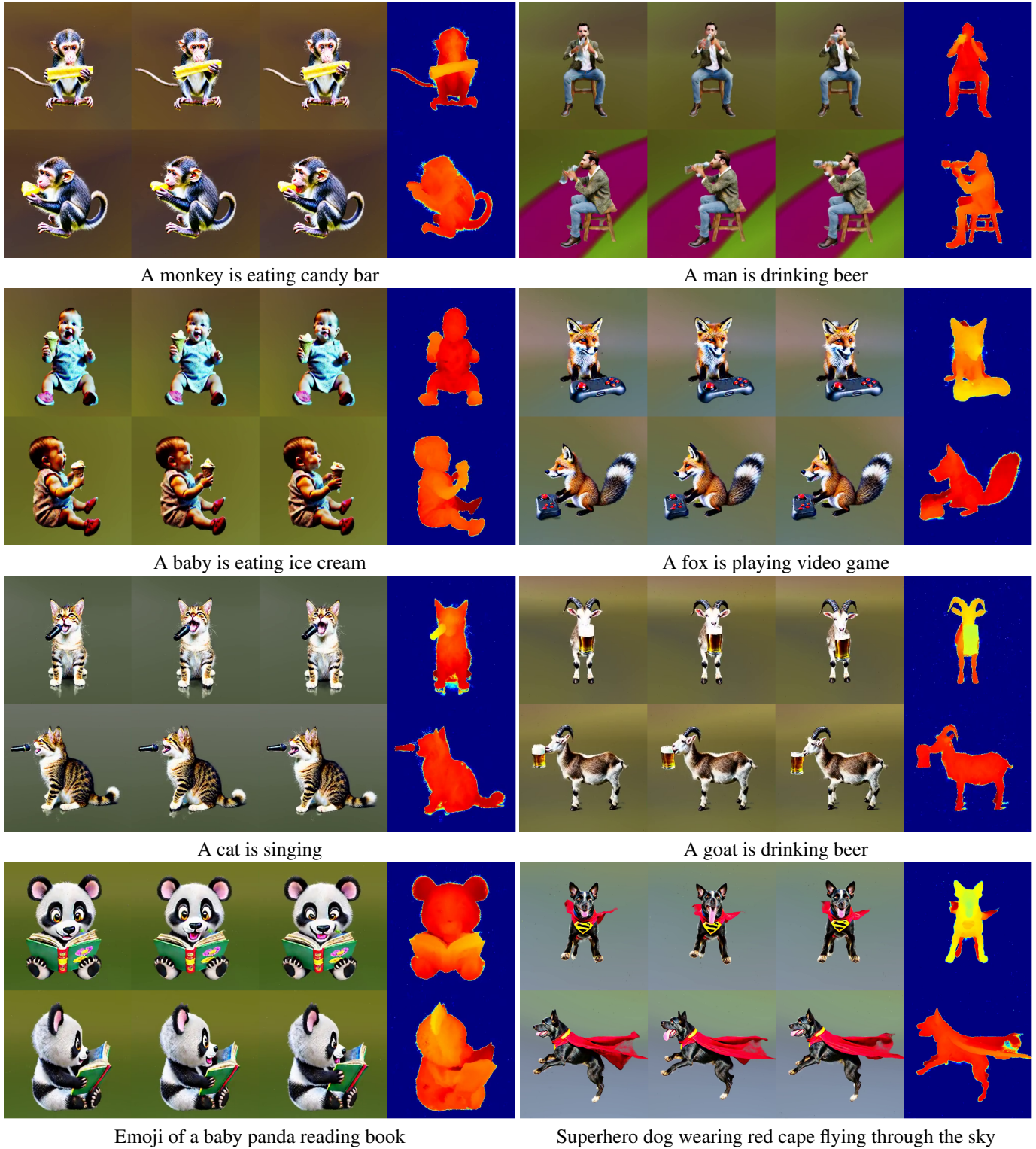Superhero dog wearing red cape flying through the sky

Figure 7. **Text-to-4D generation.** We show qualitative results of the text-to-4D generation task, demonstrating high visual quality, multi-view consistency, plausible composition and realistic motion. Video results are available in the Supplement.
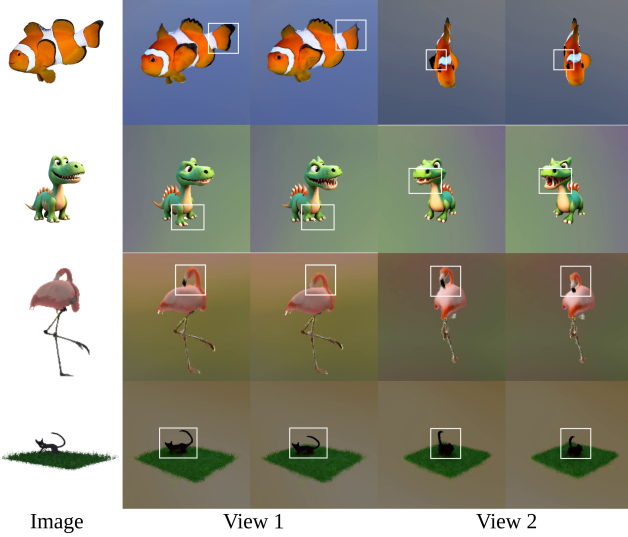
Figure 8. **Image-to-4D generation.** Given an input image, our method reconstructs and animates 3D assets. Prompts used for motion are 'Clown fish swimming', 'Cartoon dragon running', 'Flamingo scratching its neck', and 'Cat walking on the grass'. Video results can be found in the Supplement.

representation [5] used in a previous work [34]. Due to the entanglement of geometry and motion in hexplanes, it is not trivial to keep the static geometry parts frozen during the motion learning stage, which leads to lower visual quality and reappearance of the Janus problem (See Fig. 5 and Tab. 1c, 'Hexplane' versus Ours). In comparison, our dynamic representation fully disentangles the canonical model and the deformation field, successfully preserving the 3D shape and texture of the static model while learning motion.

### 4.2. Controllable 4D Generation

In this section, we show qualitative results of text-to-4D generation where the object appearance is defined by one or multiple user-defined images. In Fig. 8, we present results for single image to 4D generation. Our method can preserve the identity and appearance details of the image and successfully learn the animation specified in the text prompt. In Fig. 9, we show personalized 4D generation. Given a few casually captured images of a subject, our method can generate 4D content of the subject under various motion conditions, e.g., eating food or ice cream.

### 5. Discussion

**Conclusion** We propose Dream-in-4D, a unified approach for 4D scene synthesis from a text prompt and optionally one or multiple images as the appearance condition. Leveraging 3D and 2D diffusion priors, our method first learns a high-quality static asset, offering a good starting point for deformation optimization. Then, our motion-disentangled 4D representation allows us to learn
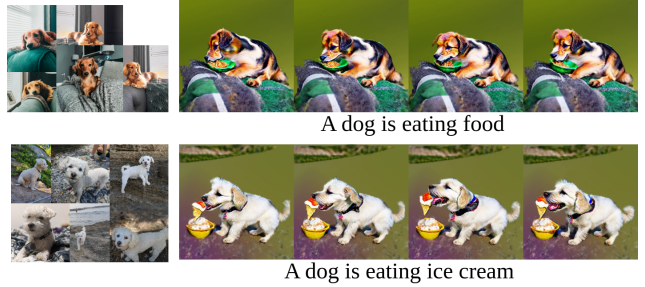


A dog is eating food



A dog is eating ice cream

Figure 9. **Personalized 4D.** Our method can generate dynamic 3D scenes of a subject given a text prompt and 4-6 causally captured images of the subject. Videos are available in the Supplement.



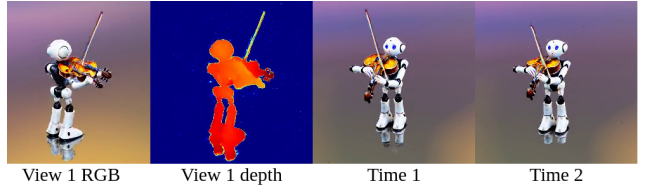View 1 RGB    View 1 depth    Time 1    Time 2

Figure 10. **Failure case.** Despite combining 3D and 2D diffusion guidance, our method fails to reconstruct some surreal prompts, e.g., 'Robot is playing the violin'. The static stage fails to learn a view-consistent violin, and the robot's hand position is incorrect. In the second stage, our method cannot correct such errors or learn plausible hand and arm motion.

motion with video diffusion guidance while maintaining the quality of the static asset. We introduce multi-resolution feature grids and a TV loss for the deformation field, resulting in more realistic motion. Dream-in-4D achieves better visual quality, 3D consistency, motion and spatial layout on the text-to-4D task compared to several baselines, while also enabling image-to-4D and personalized 4D generation.

**Limitation** Although the combination of 3D and 2D diffusion priors solves the Janus problem and achieves plausible spatial layout in most cases, it fails for some surreal prompts (e.g. robot playing violin in Fig. 10). In the dynamic stage, our method cannot recover from the wrong static representation and fails to learn correct motion given the wrong position of the hands. We believe this problem could potentially be solved with further advances in 3D and 2D diffusion models.

**Societal impact** We note that our work could be potentially used to generate fake 4D content that is violent or harmful. Building upon pre-trained large-scale diffusion models, it inherits the biases and limitation of these models. Therefore, the 4D videos generated with our method should be carefully examined and labeled as synthetic content.

## Supplemental Materials

## 1. Video Results

Please refer to our web page (https://research.nvidia.com/labs/nxp/dream-in-4d/) for video results of the text-to-4D, image-to-4D and personalized 4D tasks. We also provide qualitative comparisons of our method and ablation baselines. A low resolution version of our web page is available in *webpage_small/index.html*.

## 2. Network Architecture

**Canonical NeRF.** We use a hash-encoded multi-resolution feature grid with 16 resolution levels, where the base resolution is $16 \times 16 \times 16$ and maximum resolution is $4096 \times 4096 \times 4096$. The feature grid is followed by two shallow MLPs to produce the density and color values. Both of the MLP networks have 1 hidden layer and 64 neurons per layer.

**Deformation field.** The deformation field uses a hash-encoded multi-resolution feature grid with 12 levels, where the base resolution is $4 \times 4 \times 4 \times 4$ and maximum resolution is $232 \times 232 \times 232 \times 232$. The feature grid is followed by an MLP with 4 hidden layers and 64 neurons per layer to predict the displacement values $\mathbf{d}$ for scene deformation. We then calculate the canonical point location by $\mathbf{x_c} = \mathbf{x_d} + \mathbf{d}$.

**Background.** We model the background with an MLP, which takes the viewing direction as input and outputs a color value. This assumes that the background is located infinitely far away from the camera. The MLP has 3 hidden layers and 64 neurons per layer.

## 3. Training Schedule

### 3.1. Static Stage

For the static stage, we render multi-view images of resolution $64 \times 64$ with a batch size of 8 for the first 5000 iterations, and resolution $256 \times 256$ with a batch size of 4 for the last 5000 iterations. For MVDream [32] guidance, the images are upsampled to $256 \times 256$. For StableDiffusion [29], we upsample to $512 \times 512$. We use guidance scale of 50 for MVDream and 100 for StableDiffusion.

We use an AdamW [18] optimizer with learning rate 0.001 for all the MLP parameters and 0.01 for the parameters of hash-encoded multi-resolution feature grid. The $\beta$ parameters are set to 0.9 and 0.99. We train the networks for 10000 iterations on a NVIDIA V100 GPU, taking 4.5 hours.

_____

Websites are best viewed using Chrome, Safari or Microsoft Edge.

| Prompts | $\lambda_{2D}$ | $\lambda_{3D}$ |
|---|---|---|
| Superhero dog wearing a red cape is flying through the sky | 1.2 | 1 |
| A cat singing | 1.2 | 1 |
| A dog riding a skateboard | 1 | 1 |
| Clown fish swimming through coral reef | 1 | 1 |
| A fox playing a video game | 1.2 | 1 |
| A goat drinking beer | 1 | 1 |
| A monkey eating a candy bar | 1 | 1 |
| An emoji of a baby panda reading a book | 1.2 | 1 |
| A baby panda eating ice cream | 1 | 1 |
| A squirrel riding a motorcycle | 1.2 | 1 |

Table 1. 2D and 3D guidance loss weights for the static stage.

We found balancing the 3D and 2D guidance weight to be important for achieving view-consistent, text-aligned and realistic results. In Tab. 1, we list the loss weights for the prompts in the paper.

### 3.2. Dynamic Stage

To learn deformation in the dynamic stage, we use zero-scope guidance [3] for the main paper, where we render 24-frame videos of resolution $144 \times 80$. We found our method to also work well with modelscope guidance [2], in which case we rendered videos of $64 \times 64$ for the first 7000 iterations and then upsample to $256 \times 256$. We use a guidance scale of 100 for both zeroscope and modelscope guidance, and gradually decrease the time step used for the SDS loss [25] from $[0.99, 0.99]$ to $[0.2, 0.5]$.

In the dynamic stage, we optimize the deformation parameters with an AdamW [18] optimizer with learning rate 0.001 and $\beta = [0.9, 0.99]$. We train the deformation network for 10000 iterations on a NVIDIA A100 or RTX A6000 GPU. We start by using the first 4 levels of the multi-resolution features, and gradually include the higher resolution features, adding 1 level every 500 iterations. The dynamic stage takes 9 hours for zeroscope, and 6 hours for modelscope due to the lower inference resolution of the diffusion model.

## 4. User Study

In this section, we introduce more details of the user study. As discussed in Sec. 4.1 and Tab. 1 in the main paper, we carry out ablation studies to verify the key elements of the proposed method. Specifically, we present the results from ablated models to a user and give instructions "Please find the method that best match the prompt / has best motion / has best 3D consistency and visual quality." The example interface is shown in Fig. 1. For each user, we show videos synthesized from 10 prompts for each metric evaluation, leading to 30 questions per user in total.
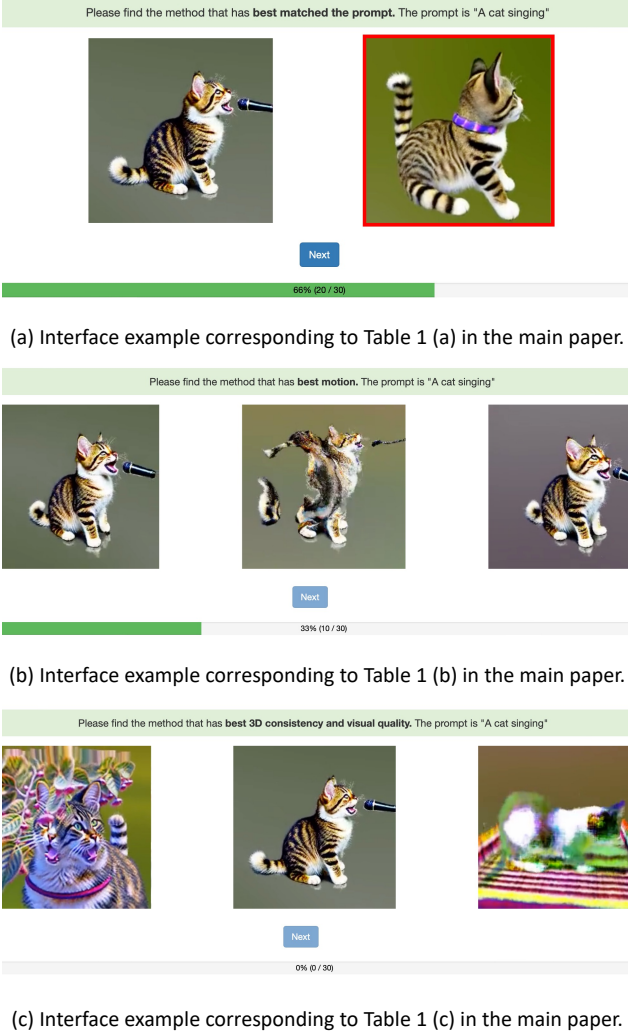
(a) Interface example corresponding to Table 1 (a) in the main paper.



(b) Interface example corresponding to Table 1 (b) in the main paper.



(c) Interface example corresponding to Table 1 (c) in the main paper.

Figure 1. **User study interface.**

# References

[1] Deepfloyd if. https://github.com/deep-floyd/IF. 1, 2, 3, 5

[2] Modelscope. https://huggingface.co/damo-vilab/text-to-video-ms-1.7b. 1, 3, 5

[3] Zeroscope. https://huggingface.co/cerspense/zeroscope_v2_576w. 2, 5, 1

[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 *CVPR*, 2023. 1, 3

[5] Ang Cao and Justin Johnson. Hexplane: A fast representation for dynamic scenes. *CVPR*, 2023. 2, 3, 8

[6] Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J Guibas, Jonathan Tremblay, Sameh Khamis, et al. Efficient geometry-aware 3d generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16123–16133, 2022. 3

[7] Rui Chen, Yongwei Chen, Ningxin Jiao, and Kui Jia. Fantasia3d: Disentangling geometry and appearance for high-quality text-to-3d content creation. In *ICCV*, 2023. 3

[8] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3d objects. *arXiv preprint arXiv:2212.08051*, 2022. 3

[9] Matt Deitke, Ruoshi Liu, Matthew Wallingford, Huong Ngo, Oscar Michel, Aditya Kusupati, Alan Fan, Christian Laforte, Vikram Voleti, Samir Yitzhak Gadre, Eli VanderBilt, Aniruddha Kembhavi, Carl Vondrick, Georgia Gkioxari, Kiana Ehsani, Ludwig Schmidt, and Ali Farhadi. Objaverse-xl: A universe of 10m+ 3d objects. *arXiv preprint arXiv:2307.05663*, 2023. 3

[10] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3d objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13142–13153, 2023. 4

[11] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5501–5510, 2022. 3

[12] Sara Fridovich-Keil, Giacomo Meanti, Frederik Rahbæk Warburg, Benjamin Recht, and Angjoo Kanazawa. K-planes: Explicit radiance fields in space, time, and appearance. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12479–12488, 2023. 3

[13] Yuan-Chen Guo, Ying-Tian Liu, Ruizhi Shao, Christian Laforte, Vikram Voleti, Guan Luo, Chia-Hao Chen, Zi-Xin Zou, Chen Wang, Yan-Pei Cao, and Song-Hai Zhang. threestudio: A unified framework for 3d content generation. https://github.com/threestudio-project/threestudio, 2023. 3

[14] Tianye Li, Mira Slavcheva, Michael Zollhoefer, Simon Green, Christoph Lassner, Changil Kim, Tanner Schmidt, Steven Lovegrove, Michael Goesele, Richard Newcombe, et al. Neural 3d video synthesis from multi-view video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5521–5531, 2022. 3

[15] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3d: High-resolution text-to-3d content creation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 300–309, 2023. 1, 3

[16] Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick. Zero-1-to-3: Zero-shot one image to 3d object, 2023. 1, 2, 3, 4, 5, 6

[17] Yuan Liu, Cheng Lin, Zijiao Zeng, Xiaoxiao Long, Lingjie Liu, Taku Komura, and Wenping Wang. Syncdreamer: Learning to generate multiview-consistent images from a

single-view image. *arXiv preprint arXiv:2309.03453*, 2023. 1, 3, 4

[18] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 1

[19] Ricardo Martin-Brualla, Noha Radwan, Mehdi SM Sajjadi, Jonathan T Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7210–7219, 2021. 3

[20] Luke Melas-Kyriazi, Christian Rupprecht, Iro Laina, and Andrea Vedaldi. Realfusion: 360 reconstruction of any object from a single image. In *CVPR*, 2023. 3

[21] Gal Metzer, Elad Richardson, Or Patashnik, Raja Giryes, and Daniel Cohen-Or. Latent-nerf for shape-guided generation of 3d shapes and textures. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12663–12673, 2023. 1, 3

[22] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Commun. ACM*, 65(1):99–106, 2021. 2, 3, 6

[23] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)*, 41(4):1–15, 2022. 3

[24] Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5865–5874, 2021. 3

[25] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022. 1, 3, 4

[26] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-NeRF: Neural Radiance Fields for Dynamic Scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020. 2, 3

[27] Guocheng Qian, Jinjie Mai, Abdullah Hamdi, Jian Ren, Aliaksandr Siarohin, Bing Li, Hsin-Ying Lee, Ivan Skorokhodov, Peter Wonka, Sergey Tulyakov, et al. Magic123: One image to high-quality 3d object generation using both 2d and 3d diffusion priors. *arXiv preprint arXiv:2306.17843*, 2023. 1, 3, 4, 6

[28] Amit Raj, Srinivas Kaza, Ben Poole, Michael Niemeyer, Ben Mildenhall, Nataniel Ruiz, Shiran Zada, Kfir Aberman, Michael Rubenstein, Jonathan Barron, Yuanzhen Li, and Varun Jampani. Dreambooth3d: Subject-driven text-to-3d generation. *ICCV*, 2023. 3

[29] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2021. 1, 2, 3

[30] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22500–22510, 2023. 6

[31] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35:36479–36494, 2022. 1

[32] Yichun Shi, Peng Wang, Jianglong Ye, Long Mai, Kejie Li, and Xiao Yang. Mvdream: Multi-view diffusion for 3d generation. *arXiv:2308.16512*, 2023. 1, 2, 3, 4

[33] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, et al. Make-a-video: Text-to-video generation without text-video data. *arXiv preprint arXiv:2209.14792*, 2022. 1, 2

[34] Uriel Singer, Shelly Sheynin, Adam Polyak, Oron Ashual, Iurii Makarov, Filippos Kokkinos, Naman Goyal, Andrea Vedaldi, Devi Parikh, Justin Johnson, and Yaniv Taigman. Text-to-4d dynamic scene generation, 2023. 2, 3, 6, 8

[35] Junshu Tang, Tengfei Wang, Bo Zhang, Ting Zhang, Ran Yi, Lizhuang Ma, and Dong Chen. Make-it-3d: High-fidelity 3d creation from a single image with diffusion prior, 2023. 3

[36] Christina Tsalicoglou, Fabian Manhardt, Alessio Tonioni, Michael Niemeyer, and Federico Tombari. Textmesh: Generation of realistic 3d meshes from text prompts. *arXiv preprint arXiv:2304.12439*, 2023. 3

[37] Haochen Wang, Xiaodan Du, Jiahao Li, Raymond A. Yeh, and Greg Shakhnarovich. Score jacobian chaining: Lifting pretrained 2d diffusion models for 3d generation. *arXiv preprint arXiv:2212.00774*, 2022. 3

[38] Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation. *arXiv preprint arXiv:2305.16213*, 2023. 3

[39] Joseph Zhu and Peiye Zhuang. Hifa: High-fidelity text-to-3d with advanced diffusion guidance. *arXiv preprint arXiv:2305.18766*, 2023. 3, 5