# STABLE VIRTUAL CAMERA: Generative View Synthesis with Diffusion Models

Jensen (Jinghao) Zhou[1,2,*,†]    Hang Gao[1,3,*,†]

Vikram Voleti[1]    Aaryaman Vasishta[1]    Chun-Han Yao[1]    Mark Boss[1]

Philip Torr[2]    Christian Rupprecht[2]    Varun Jampani[1]

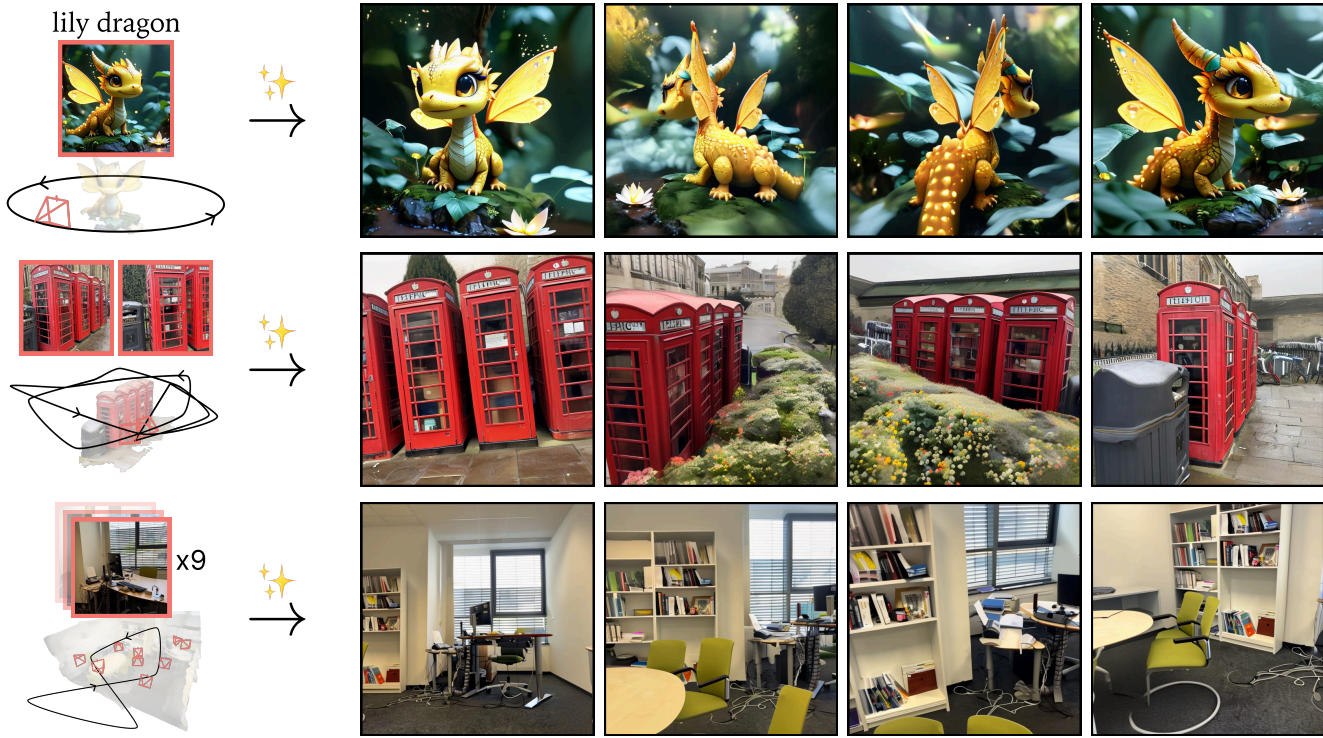[1]Stability AI    [2]University of Oxford    [3]University of California, Berkeley

Figure 1. **Generative view synthesis.** STABLE VIRTUAL CAMERA generates novel views from any number of input views and target cameras, which the user can specify anywhere. We show three examples: single view with simple orbit camera trajectory (top); two views with long camera trajectory (middle); and nine views with large spatial range (bottom). Please visit our website for video samples.

## Abstract

*We present STABLE VIRTUAL CAMERA (SEVA), a generalist diffusion model that creates novel views of a scene, given any number of input views and target cameras. Existing works struggle to generate either large viewpoint changes or temporally smooth samples, while relying on specific task configurations. Our approach overcomes these limitations through simple model design, optimized training recipe, and flexible sampling strategy that generalize across view synthesis tasks at test time. As a result, our samples maintain high consistency without requiring additional 3D representation-based distillation, thus streamlining view synthesis in the wild. Furthermore, we show that our method can generate high-quality videos lasting up to half a minute with seamless loop closure. Extensive benchmarking demonstrates that SEVA outperforms existing methods across different datasets and settings.*

---

*Equal contribution.

†Jensen and Hang did the work during internships at Stability AI.

# 1. Introduction

Novel view synthesis (NVS) aims to generate realistic, 3D-consistent images of a scene from arbitrary camera viewpoints given any number of camera-posed input views. Traditional methods, which rely on dense input views, treat NVS as a 3D reconstruction and rendering problem [1–3], but this approach fails with sparse inputs. Generative view synthesis addresses this limitation by leveraging modern deep network priors [4, 5], enabling immersive 3D interactions in uncontrolled environments without the need to capture large image sets per scene. In this work, we focus on generative view synthesis and, unless otherwise specified, refer to it simply as NVS for clarity.

Despite recent progress [6–12], NVS in the wild remains limited due to two key challenges: First, existing methods struggle to generate both large viewpoint changes [9, 10] and temporally smooth samples [6–8, 13] while being constrained by rigid task configurations, such as a fixed number of input and target views [7, 9, 11, 12], reviewed in Tab. 1. Second, their sampling consistency is often insufficient, necessitating additional NeRF distillation to fuse inconsistent results into a coherent representation [7, 8, 13]. These limitations hinder their applicability across diverse NVS tasks, which we address in this work.

We present S̲TABLE V̲IRTUAL C̲AMERA[1] (SEVA), a diffusion-based NVS model that generalizes across a spectrum of view synthesis tasks without requiring NeRF distillation. With a single network, SEVA generates high-quality novel views that strike both large viewpoint changes and temporal smoothness, while supporting any number of input and target views. Our approach simplifies the NVS pipeline without requiring distillation from a 3D representation, thus streamlining it for real-world applications. For the first time, we demonstrate high-quality videos lasting up to half a minute with precise camera control and seamless loop closure in 3D. We highlight these results in Fig. 1 and showcase more examples of camera control in Fig. 2.

To achieve this, we carefully design our pipeline in three key aspects: model design, training recipe, and sampling method at inference. First, SEVA avoids explicit 3D representations within the network, allowing the model to inherit strong priors from pre-trained 2D models. Second, during training, we carefully craft our view selection strategy to cover both small and large viewpoint changes, ensuring strong generalization to diverse NVS tasks. Third, at inference, we introduce a two-pass procedural sampling approach that supports flexible input-target configurations. Together, these design choices create a versatile 3D "virtual camera simulation system" capable of synthesizing novel views along arbitrary camera trajectories with any number

---

[1]We name this model in tribute to the Virtual Camera [14] cinematography technology, a pre-visualization technique to simulate real-world camera movements.

| model | training data | generation capacity | interpolation smoothness | input flexibility |
|---|---|---|---|---|
| **Regression-based** | | | | |
| pixelNeRF [4] | ◉ | ✗ | ✓ | sparse (1) |
| pixelSplat [16] | ▲ | ✗ | ✓ | sparse (2) |
| MVSplat [17] | ▲ | ✗ | ✓ | sparse (2) |
| Long-LRM [18] | ▲ | ✗ | ✓ | semi-dense ({16, 32}) |
| LVSM [11] | ◉/▲ | ✗ | ✓ | sparse ({2, 4}) |
| **Diffusion-based: image models** | | | | |
| Zero123 [6] | ◉ | ✓ | ✗ | sparse (1) |
| ZeroNVS [13] | ◉▲ | ✓ | ✗ | sparse (1) |
| ReconFusion [7] | ◉▲ | ✓ | ✗ | sparse (3) |
| CAT3D [8] | ◉▲ | ✓ | ✗ | sparse ([1, 9]) |
| **Diffusion-based: video models** | | | | |
| SV3D [19] | ◉ | ✗ | ✓ | sparse (1) |
| MotionCtrl [10] | ▲ | ✗ | ✓ | sparse (1) |
| ViewCrafter [9] | ▲ | ✗ | ✓ | sparse (2) |
| 4DiM [12] | ▲ | ✓ | ✓ | sparse ({1, 2, 8}) |
| SEVA | ◉▲ | ✓ | ✓ | sparse ([1, 8]), semi-dense ([9, 32*]) |

Table 1. **Comparison of existing NVS models** based on the source of training data and key attributes. SEVA is trained on both object-level (◉) and scene-level (▲) data, offering flexible input conditioning, strong generation capacity, and smooth view interpolation. We define generation capacity and interpolation smoothness of each work based on their evaluation setting and our benchmark results. *This upper-bound can be up to hundreds for dense captures, we test our model up to 32 views in practice.

of input and target views, without using a 3D representation.

We conducted a unified benchmark across 10 datasets and a variety of experimental settings, including both open-source and proprietary models. Our benchmark reflects the diversity of real-world NVS tasks across the board and systematically evaluates existing methods beyond their comfort zones. We find that SEVA consistently outperforms previous works, achieving +1.5 dB PSNR over the state of the art CAT3D [8] in its own setup. Moreover, our method generalizes well to in-the-wild user captures, with input views ranging from 1 to 32.

In summary, our key contributions with the SEVA model include: (1) a training strategy for jointly modeling large viewpoint changes and temporal smoothness, (2) a two-pass procedural sampling method for smooth video generation along arbitrary long camera trajectories, (3) a comprehensive benchmark that evaluates NVS methods across different datasets and settings, and (4) an open-source release of model weights to support future research.

## 2. Background

We consider the evaluation of an NVS model across three key criteria: (1) generation capacity—the ability to synthesize missing regions for large viewpoint changes; (2) interpolation smoothness—the ability to produce seamless transitions between views; and (3) input flexibility—the abil-
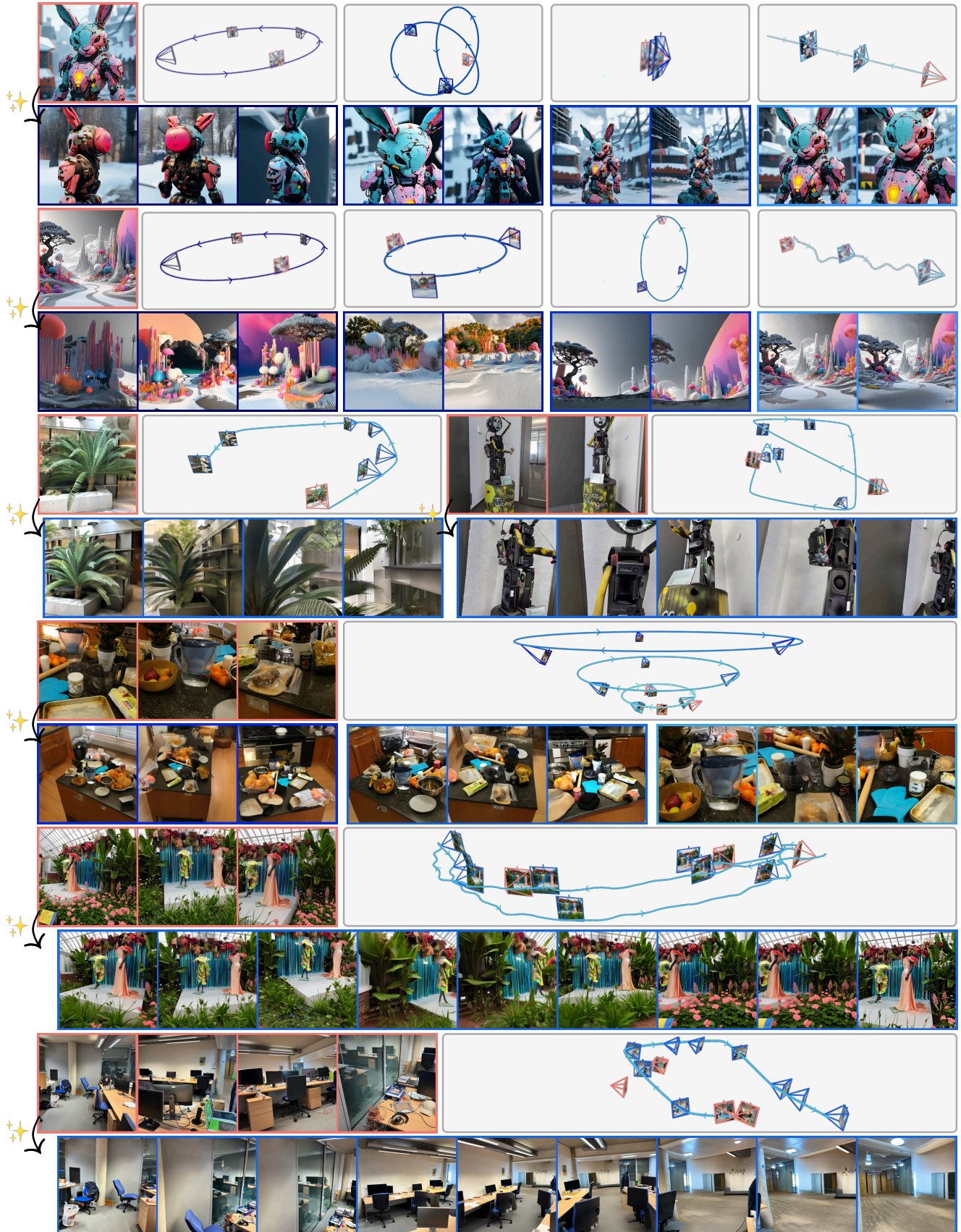
Figure 2. **Diverse camera control.** SEVA generates photorealistic novel views following diverse camera trajectories. This includes orbit, spiral, zoom out, dolly zooms, and any user-specified trajectories. Please visit our website for more visual results.

ity to handle a variable number of input and target views; We review existing NVS models based on these criteria in Tab. 1, including the types of training data.

## 2.1. Types of NVS Tasks

Given $M$ input view images $\mathbf{I}^{\text{inp}} \in \mathbb{R}^{M \times H \times W \times 3}$ of $H \times W$ resolution, along with their corresponding cameras $\boldsymbol{\pi}^{\text{inp}}$, NVS involves predicting $N$ targets views $\mathbf{I}^{\text{tgt}} \in \mathbb{R}^{N \times H \times W \times 3}$, specified by their respective cameras $\boldsymbol{\pi}^{\text{tgt}}$. For each camera, we assume we know both intrinsics and extrinsics. Based on the number of input views, we define the "sparse-view regime" as having up to 8 input views, and the "semi-dense-view regime" as an intermediate state bridging the sparse-view regime and dense captures, which typically involve hundreds of views. Based on the nature of their target views, we bucket a broad range of NVS tasks into "set NVS" and "trajectory NVS", as shown in Fig. 3.

**Set NVS** considers a set of target views in arbitrary order, usually across a large spatial range. The order of views is often not helpful here, and a good NVS model requires great generation capacity to excel at this task. We note that some works address only this task (*e.g.* ReconFusion [7]).

**Trajectory NVS** regards target views along a smooth camera trajectory, such that they form a video sequence. However, they are often sampled within a small spatial range in a shorter video. To solve this task, a good NVS model requires great interpolation smoothness to produce consistent and non-flickering results. We note that some existing works address only this task (*e.g.* ViewCrafter [9]).

## 2.2. Existing Models

We group existing approaches into regression- and diffusion-based models based on their high-level design choices. A more detailed discussion of related works can be found in Appendix B.

**Regression-based models** learn a deterministic mapping:

$$f_\theta\left(\mathbf{I}^{\text{inp}}, \boldsymbol{\pi}^{\text{inp}}, \boldsymbol{\pi}^{\text{tgt}}\right),$$

to directly generate $\mathbf{I}^{\text{tgt}}$ deterministically from $\mathbf{I}^{\text{inp}}, \boldsymbol{\pi}^{\text{inp}}, \boldsymbol{\pi}^{\text{tgt}}$. $f_\theta$ can be either an end-to-end network parameterized by $\theta$, or a composition of a feed-forward prediction of an intermediate 3D representation and then a neural renderer (*e.g.*, NeRF [20] or 3DGS [3]). For the latter case, set NVS and trajectory NVS are solved in the same way since there exists a persistent 3D representation.

**Diffusion-based models** capture the conditional distribution:

$$p_\theta\left(\mathbf{I}^{\text{tgt}} \mid \mathbf{I}^{\text{inp}}, \boldsymbol{\pi}^{\text{inp}}, \boldsymbol{\pi}^{\text{tgt}}\right),$$

from which $\mathbf{I}^{\text{tgt}}$ are sampled [21] iteratively. We highlight two types of models within this scope: Image and



Figure 3. **Set NVS *versus* trajectory NVS.** Set NVS generates target views as an image set, whereas trajectory NVS produces them as a trajectory video.

Video models. Image models are trained on unordered image sets, such that $(\mathbf{I}^{\text{inp}}, \mathbf{I}^{\text{tgt}}) \sim \mathcal{I}$, where $\mathcal{I} = \{\mathbf{I}_{\sigma(1)}, \mathbf{I}_{\sigma(2)}, \cdots, \mathbf{I}_{\sigma(M+N)}\}$ is an image batch, and $\sigma(\cdot)$ is a random permutation function, where camera parameters are omitted for simplicity. Image models usually thrive at set NVS, but struggle in trajectory NVS since they are designed to generate images and not videos. Additionally, the unordered nature of all views solicits flexible input conditioning. Video models are instead trained on ordered views, such that $(\mathbf{I}^{\text{inp}}, \mathbf{I}^{\text{tgt}}) \sim \mathcal{V}$, where $\mathcal{V} = \{\mathbf{I}_1, \mathbf{I}_2, \cdots, \mathbf{I}_{M+N}\}$ is a randomly sampled video batch with ordering preserved. Additional temporal operators may also be used to improve the temporal smoothness, such as temporal positional encoding and temporal attention. In contrast with image models, video models thrive at trajectory NVS, but struggle in set NVS. Moreover, all existing video models require both input and target views to be ordered (input views followed by target ones), constraining their input flexibility [10, 19, 22–24].

## 2.3. Remarks and Motivation

Existing tasks pose critical challenges to our design choices. Specifically, our design choices are made to achieve high generation capacity, smooth view interpolation, and flexible input conditioning, as compared in Tab. 1. In this way, we can employ a single model for both tasks, described next.

## 3. Method

We describe our model design and training strategy in Sec. 3.1, then our sampling process at test time in Secs. 3.2 and 3.3. A system overview is provided in Fig. 4.

## 3.1. Model Design and Training

We consider a "$M$-in $N$-out" multi-view diffusion model $p_\theta$, as notated in Sec. 2.2. We formulate this learning problem as a standard diffusion process [21] without any change.

**Architecture.** Our model is based on the publicly available SD 2.1 [25], which consists of an auto-encoder and a latent denoising U-Net. Following [8], we inflate the 2D self-attention of each low-resolution residual block into 3D self-attention [26] within the U-Net. To improve model capacity, we add 1D self-attention along the view axis after each self-attention block via skip connection [27, 28],
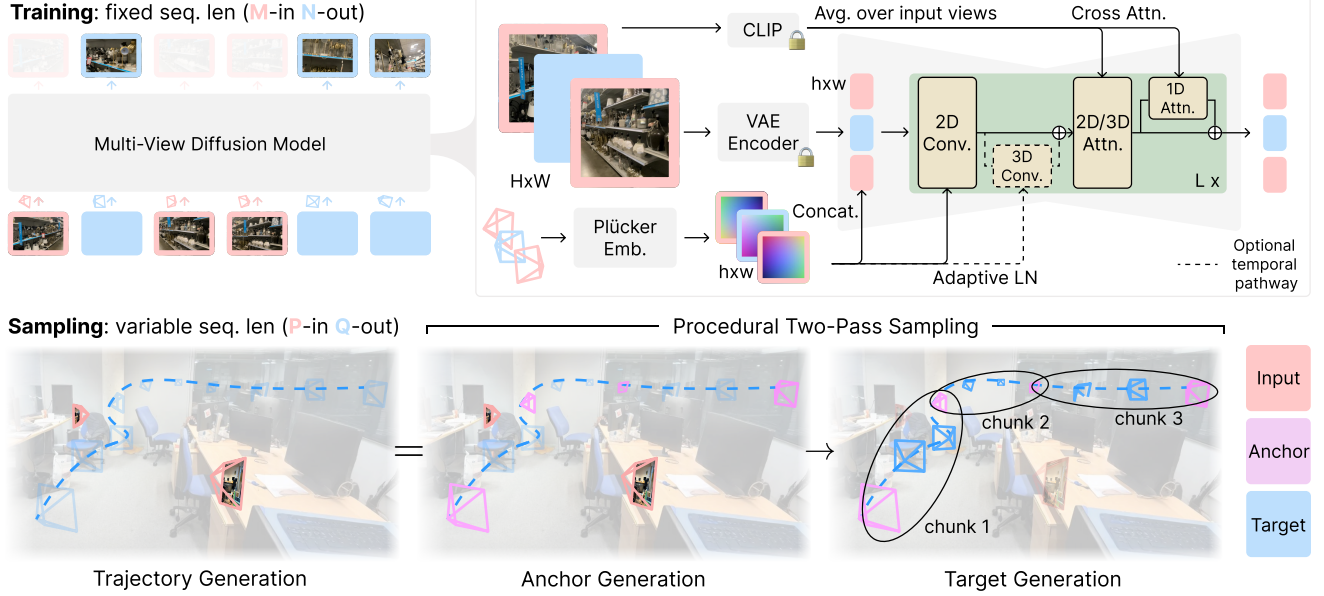
Figure 4. **Method.** SEVA is trained with fixed sequence length as a "$M$-in $N$-out" multi-view diffusion model with standard architecture. It conditions on CLIP embeddings, VAE latents of the input views, and their corresponding camera poses. During sampling, SEVA can be cast as a generative "$P$-in $Q$-out" renderer that works with variable sequence length, where $P$ and $Q$ need not be equal to $M$ and $N$. To enhance temporal and 3D consistency across generated views, especially when generating along a trajectory, we present procedural two-pass sampling as a general strategy.

bumping the model parameters from 870M to 1.3B. Optionally, we further tame this model into a video model by introducing 3D convolutions after each residual block via skip connection, similar to [22, 29], yielding 1.5B total parameters. The temporal pathway can be enabled during inference when frames within one forward pass are known to be spatially ordered, enhancing output's smoothness.

**Conditioning.** To fine-tune our base model into a multi-view diffusion model, we add camera conditioning as Plücker embedding [30] via concatenation [8] and adaptive layer normalization [31]. We normalize $\pi^{\text{inp}}$ and $\pi^{\text{tgt}}$ by first computing the relative pose with respect to the first input camera and then normalizing the scene scale such that all camera positions are within a $[-2, 2]^3$ cube. For each input frame, we first encode its latent then concatenate with its Plücker embedding and a binary mask [8, 22] differentiating between input and target views. For each target frame, we use the noisy state of its latent instead. Additionally, we find it helpful [19] to also inject high-level semantic information via CLIP [32] image embedding. We zero initialize new weights for additional channels in the first layer. In our experiment, we found that our model can quickly adapt to these conditioning changes and produce realistic images with as few as $5K$ iterations.

**Training.** Let us define the training context window length $T = |\mathbf{I}^{\text{inp}}| + |\mathbf{I}^{\text{tgt}}| = M + N$. One natural goal is to support large $T$ such that we can generate a larger set of frames. However, we find that naive training is prone to

divergence, and we thus employ a two-stage training curriculum. During the first stage, we train our model with $T = 8$ with a batch size of $1472$ for $100K$ iterations. In the second stage, we train our model with $T = 21$ with a batch size of $512$ for $600K$ iterations. Given a training video sequence, we randomly sample the number of input frames $M \in [1, T-1]$ and the frames $(\mathbf{I}^{\text{inp}}, \mathbf{I}^{\text{tgt}})$. We find it important to jointly sample $\mathbf{I}$ with a smaller subsampling stride to ensure sufficient temporal granularity and avoid missing critical transitions with a small probability ($0.2$ is used in practice). In the optional video training stage, we only train temporal weights with data sampled with a small subsampling stride and a batch size of $512$ for $200K$ iterations. We shift the signal-to-noise ratio (SNR) in all stages as more noise is necessary to destroy the information when training with more frames, corroborating findings from [8, 33, 34]. The model is trained with squared images with $H = W = 576$.

### 3.2. Sampling Novel Views

Once the diffusion model is trained, we can sample it for a wide range of NVS tasks during test time. Formally, let us consider a "$P$-in $Q$-out" NVS task during testing, where we are given $P = |\mathbf{I}^{\text{inp}}|$ input frames and aim to produce $Q = |\mathbf{I}^{\text{tgt}}|$ target frames. Our goal is to design a generic sampling strategy that works for all $P$ and $Q$ configurations, where $P$ and $Q$ need not be equal to $M$ and $N$.

We make two key observations: First, within a single forward pass, predictions are 3D consistent, provided

the model is well-trained. Second, when $P + Q > T$, $\mathbf{I}^{\text{tgt}}$ must be split into smaller chunks of size $Q_i$ such that $P + Q_i \leqslant T$ for the $i^{\text{th}}$ forward pass. We term this practice *one-pass sampling*. However, predictions across these forward passes would be inconsistent unless they share common frames to maintain local consistency within a spatial neighborhood. Building on these observations, we summarize our sampling process under two scenarios: $P + Q \leqslant T$ and $P + Q > T$.

**$P + Q \leqslant T$.** We fit the task within one forward pass for simplicity and consistency. As shown in Appendix D, we find it works better to pad the forward pass to have exactly $T$ frames by repeating the first input image, compared to changing the context window $T$ zero-shot.

**$P + Q > T$.** We propose *procedural two-pass sampling*: In the first pass, we generate anchor frames $\mathbf{I}^{\text{acr}}$ using all input frames $\mathbf{I}^{\text{inp}}$. In the second pass, we divide $\mathbf{I}^{\text{tgt}}$ into chunks and generate them using $\mathbf{I}^{\text{acr}}$ (and optionally $\mathbf{I}^{\text{inp}}$) according to the spatial distribution of $\mathbf{I}^{\text{acr}}$ and $\mathbf{I}^{\text{tgt}}$. Given the distinct nature of the two tasks of interest—set NVS and trajectory NVS—*e.g.*, differences in the availability of views' ordering, we design tailored chunking strategies for each task.

For set NVS, we consider $\text{nearest}$ procedural sampling. We first generate $\mathbf{I}^{\text{acr}}$ based on pre-defined trajectory priors, similar to [8], *e.g.*, 360 trajectories for object-centric scenes, or spiral trajectories for forward-facing scenes. We then divide $\mathbf{I}^{\text{tgt}}$ into chunks w.r.t. $\mathbf{I}^{\text{acr}}$ using nearest neighbor. Specifically, the $i^{\text{th}}$ forward pass involves:

$$\text{nearest} : \{\mathbf{I}_i^{\text{acr}}\} \cup \{\mathbf{I}_j^{\text{tgt}} \mid \text{NN}(\mathbf{I}_j^{\text{tgt}}, \mathbf{I}^{\text{acr}}) = \mathbf{I}_i^{\text{acr}}\}.$$

We considered two strategies of procedural sampling: $\text{nearest}$ as described above, and $\text{gt} + \text{nearest}$ strategy by appending $\mathbf{I}^{\text{inp}}$ into each forward pass. We find that the $\text{gt} + \text{nearest}$ strategy performs better than $\text{nearest}$ and thus default to it instead. In the absence of trajectory priors, we revert to one-pass sampling. In practice, employing nearest anchors enhances qualitative consistency, albeit on a limited scale.

For trajectory NVS, we consider $\text{interp}$ procedural sampling. We first generate a subset of target frames as $\mathbf{I}^{\text{acr}}$ by uniformly spanning the target camera path with a stride $\Delta = \lfloor \frac{Q}{T-2} \rfloor$. We then generate the rest of $\mathbf{I}^{\text{tgt}}$ as segments between those anchors:

$$\text{interp} : \{\mathbf{I}_i^{\text{acr}}, \mathbf{I}_{i \cdot \Delta + 1}^{\text{tgt}}, \cdots, \mathbf{I}_{(i+1) \cdot \Delta - 1}^{\text{tgt}}, \mathbf{I}_{i+1}^{\text{acr}}\}.$$

Since the input to the model is ordered, we can leverage temporal weights to further improve smoothness (Sec. 4.3). Similarly, $\text{gt} + \text{interp}$ is possible by appending $\mathbf{I}^{\text{inp}}$ with $\Delta = \lfloor \frac{Q}{T-P-2} \rfloor$. We find that $\text{interp}$ is sufficiently robust, and choose it as the default option. The $\text{interp}$ strategy drastically outperforms its counterparts (*e.g.*, one-pass,



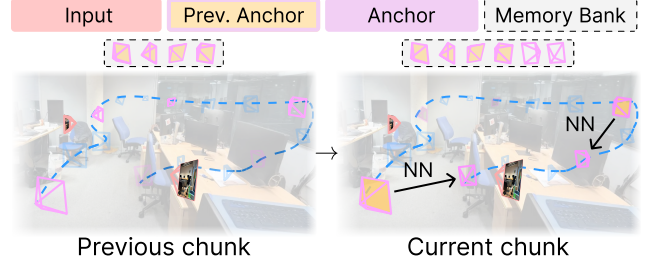Input   Prev. Anchor   Anchor   Memory Bank

Previous chunk   →   Current chunk

Figure 5. **Anchor generation when $Q >> T$.** We introduce a memory bank composed of previously-generated anchor views and their corresponding camera poses. The lookup of spatial neighbors helps improve long-term 3D consistency.

or $\text{gt} + \text{nearest}$ procedural sampling) in terms of temporal smoothness.

### 3.3. Scaling Sampling for Large $P$ and $Q$

Next, we examine two special cases when $P + Q > T$: $P > T$ and $Q \gg T$. Here, we make a tailored design for anchor generation in the first pass, while keeping target generation in the second pass unchanged.

**$P > T$.** In the semi-dense-view regime (*e.g.*, $P = 32$), we extend the context window length $T$ zero-shot to accommodate all $P$ input views and anchor views in one pass during anchor generation. Empirically, $T$ can even be extended up to hundreds without severe degradation in photorealism in the generated outputs. We find that the diffusion model generalizes well in this case as long as the input views cover the majority of the scene, shifting the task from generation to primarily interpolation. In the sparse-view regime (*i.e.*, $P \leq 8$), we observe similar performance degradation caused by zero-shot extension of $T$ compared to what we have found when $P + Q \leqslant T$. Refer to Sec. 4.5 for a detailed discussion.

**$Q \gg T$.** When the number of target views $Q$ is large, *e.g.*, in large-set NVS or long-trajectory NVS, even anchors will be chunked into different forwards in the first pass, leading to the inconsistency of anchors. To this end, we maintain a *memory bank* of anchor views previously generated, as shown in Fig. 5. We generate anchors autoregressively by retrieving their spatially nearest ones from the memory bank, similar to the $\text{nearest}$ strategy introduced above for the second pass. In Sec. 4.4, we show that this strategy drastically outperforms the standard practice of reusing temporally nearest anchors previously generated in long video literature [24], in terms of long-range 3D consistency, especially for hard trajectories.

## 4. Experiments

We employ a single model for a spectrum of settings and find that SEVA model generalizes well under the three criteria (Tab. 1). We cover different NVS tasks (set NVS and

trajectory NVS) and examine one special task of interest—long trajectory NVS. We also cover different input regimes (single-view, sparse-view, and semi-dense-view). A discussion about several key properties is presented in Sec. 4.5.

## 4.1. Benchmark

**Datasets, splits, and the number of input views.** We consider (1) object datasets, e.g., OmniObject3D [35] (OO3D) and GSO [36]; (2) object-centric scene datasets, e.g., LLFF [37], DTU [38], CO3D [39], and Wil-dRGBD [40] (WRGBD); and (3) scene datasets, e.g., RealEstate10K [41] (RE10K), Mip-NeRF 360 [42] (Mip360), DL3DV140 [43] (DL3DV), and Tanks and Temples [44] (T&T). We consider a wide range of the number of input views $P$, ranging from sparse-view regime to semi-dense-view regime, evaluating models' input flexibility. To establish a comprehensive and rigorous comparison with baselines, we consider different dataset splits utilized in prior works with the same input-view configuration, unless specified as our split (O). These include splits used in 4DiM [12] (D), ViewCrafter [9] (V), pixelSplat [16] (P), ReconFusion [7] (R), SV3D [19] (S), and Long-LRM [18] (L). For example, the 4DiM [12] (D) split on the RE10K dataset is 128 out of all 6711 test scenes with $P = 1$.

**Small-viewpoint *versus* large-viewpoint NVS.** Sweeping across all datasets, splits, and input-view configurations reveals a diverse benchmark of setups. To better evaluate models' generation capacity and interpolation smoothness (Sec. 2.1), we propose to categorize these setups into two groups—*small-viewpont* NVS and *large-viewpoint* NVS—depending on the disparity between $\mathbf{I}^{\text{tgt}}$ and $\mathbf{I}^{\text{inp}}$. small-viewpoint NVS with smaller disparities emphasizes interpolation smoothness and continuity with nearby input views, whereas large-viewpoint NVS with larger disparities requires a model to generate prominent unseen areas from input observations, predominantly assessing models' generation capacity. See Tab. 7 for the complete list. Refer to Appendix C for the detailed choice of datasets, splits, the number of input views, and the way to measure disparity.

**Baselines.** We consider a range of proprietary models, including ReconFusion [7], CAT3D [8], 4DiM [12], LVSM [11], and Long-LRM [18]. We also consider various open-source models, including SV3D [19], MVSplat [17], depthSplat [45], MotionCtrl [46], and ViewCrafter [9]. As outlined in Sec. 2.2, these baselines encompass both regression-based and diffusion-based approaches, providing a comprehensive framework for comparison.

## 4.2. Set NVS

In this section, we focus on comparing our model against prior works, given that set NVS is a task that has been ex-

tensively explored.

**Quantitative comparison.** The input and target views are chosen following splits used in previous methods. The order of target views is not preserved, *i.e.*, $\mathbf{I}^{\text{tgt}} \sim \mathcal{I}$. We use standard metrics of peak signal-to-noise ratio (PSNR), learned perceptual image patch similarity [48] (LPIPS), and structural similarity index measure [49] (SSIM). Only PSNR is showcased here due to space limits, with the rest deferred to Appendix D.2. Empirically, our method shows a greater performance improvement on LPIPS, reflecting the photo-realism of our results.

For small-viewpoint set NVS, Tab. 2 shows that SEVA sets state-of-the-art results in the majority of splits. In the sparse-view regime (*i.e.*, $P \leqslant 8$), SEVA excels across different datasets when $P > 1$. For example, a performance gain of +6.0 dB PSNR is achieved on LLFF with $P = 3$. In the semi-dense-view regime (*e.g.*, $P = 32$), SEVA surprisingly performs favorably against the specialized model [18], despite not being specifically designed for this setup. For example, SEVA lags behind the state-of-the-art method [18] by only 1.7 dB on T&T. On object datasets OO3D and GSO, SEVA achieves a significantly higher state-of-the-art PSNR compared to all other methods.

Notably, for small-viewpoint set NVS on the RealEstate10K [41] dataset, SEVA underperforms when in the single-view regime (*i.e.*, $P = 1$). This issue arises from scale ambiguity in the model due to two factors: (1) it always takes in unit-normalized cameras during training, and (2) it is trained on multiple datasets with diverse scales. This challenge is most pronounced on RE10K, where panning motion dominates. Additionally, the absence of a second input view negates any scale relativity. To address this, for all results with $P = 1$, we sweep the *unit length for camera normalization* from 0.1 to 2.0 (with 2.0 used during training), selecting the best scale for each scene. On P split with $P = 2$, we observe diffusion models lag behind regression-based models that are advantageous in small-viewpoint interpolation. SEVA bridges this gap by improving upon the state-of-the-art diffusion-based model by +4.2 dB. On R split with $P = 3$, the advantage of SEVA is pronounced exceeding the previously best result by +1.9 dB. Notably, ViewCrafter excels on V split due to capacity taking in wide-aspect-ratio images and thus more input pixels than others with square images. The advantage of ViewCrafter on V split diminishes on CO3D since the majority of informative pixels are centrally located.

For large-viewpoint set NVS, Tab. 3 shows that SEVA's quantitative advantages are even more prominent here, revealing clear benefits of SEVA in terms of generation capacity when the camera spans a large spatial range. On Mip360 with $P = 3$, SEVA improves over previous state-of-the-art method CAT3D [8] by +0.6 dB PSNR. On harder scenes like DL3DV and T&T with different input-view configura-

| Method | OO3D | GSO | RE10K | | | | | LLFF | | DTU | | CO3D | | WRGBD | | Mip360 | DL3DV | | T&T | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| split | O | O | D[12] | V[9] | P[16] | R[7] | R[7] | R[7] | R[7] | R[7] | R[7] | V[9] | R[7] | $O_e$ | $O_h$ | R[7] | O | L[18] | V[9] | L[18] |
| $P$ | 3 | 3 | 1 | 1 | 2 | 1 | 3 | 1 | 3 | 1 | 3 | 1 | 3 | 3 | 6 | 6 | 6 | 32 | 1 | 32 |
| **Regression-based models** | | | | | | | | | | | | | | | | | | | | |
| Long-LRM [18] | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 23.86 | - | 18.20 |
| MVSplat [17] | 14.78 | 15.21 | 20.42 | 20.32 | 26.39 | 21.56 | 25.64 | 11.23 | 12.50 | 13.87 | 15.52 | 12.52 | 13.52 | 14.56 | 12.54 | 13.56 | 14.34 | 16.24 | 13.22 | 12.63 |
| DepthSplat [45] | 15.67 | 16.52 | 20.90 | 19.24 | 27.44 | 21.87 | 22.54 | 12.07 | 12.62 | 14.15 | 16.24 | 13.23 | 13.77 | 15.93 | 14.23 | 14.01 | 15.72 | 16.78 | 14.35 | 13.12 |
| LVSM [11] | - | - | - | - | 29.67 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - |
| **Diffusion-based models** | | | | | | | | | | | | | | | | | | | | |
| MotionCtrl [46] | - | - | 12.74 | 16.29 | - | - | - | - | - | - | - | 15.46 | - | - | - | - | - | - | 13.29 | - |
| 4DiM [12] | - | - | 17.08 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - |
| ViewCrafter [9] | 14.64 | 15.93 | 20.43 | 22.04 | 21.42 | 20.88 | 22.81 | 10.53 | 13.52 | 12.66 | 16.40 | 18.96 | 14.72 | 16.42 | 12.66 | 14.59 | 13.78 | - | 18.07 | - |
| SEVA | 30.30 | 31.53 | 17.99 | 18.56 | 25.66 | 18.11 | 27.57 | 14.03 | 19.48 | 14.47 | 20.82 | 18.40 | 19.25 | 19.75 | 18.91 | 16.70 | 17.80 | 20.96 | 15.16 | 16.50 |

Table 2. **PSNR↑ on small-viewpoint set NVS.** $P$ denotes the number of input views. For all results with $P = 1$, we sweep the unit length for camera normalization due to the model's scale ambiguity. $O_e$ and $O_h$ denote the easy and hard split of our split, respectively. Underlined numbers are run by us using the offical released code.

| Method | OO3D | GSO | CO3D | WRGBD | | Mip360 | | DL3DV | | T&T | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| split | S[19] | S[19] | R[7] | $O_h$ | | R[7] | | O | | O | | | |
| $P$ | 1 | 1 | 1 | 1 | 3 | 1 | 3 | 1 | 3 | 1 | 3 | 6 | 9 |
| SV3D [19] | 19.28 | 20.38 | - | - | - | - | - | - | - | - | - | - | - |
| DepthSplat [45] | 11.56 | 12.32 | 10.42 | 9.35 | 13.53 | 10.49 | 12.54 | 9.63 | 12.52 | 8.63 | 9.78 | 10.12 | 11.20 |
| CAT3D [8] | - | - | - | - | - | 15.15 | - | - | - | - | - | - | - |
| ViewCrafter [9] | 10.56 | 11.42 | 10.11 | 9.12 | 13.45 | 9.79 | 10.34 | 8.97 | 11.50 | 9.23 | 9.88 | 10.32 | 11.08 |
| SEVA | 19.25 | 20.65 | 15.30 | 14.37 | 17.28 | 12.93 | 15.78 | 13.01 | 15.95 | 11.28 | 12.65 | 13.80 | 14.72 |

Table 3. **PSNR↑ on large-viewpoint set NVS.** For all results with $P = 1$, we sweep the unit length for camera normalization due to the model's scale ambiguity. Underlined numbers are run by us using the officially released code.

| Method | small-viewpoint | | | | large-viewpoint |
| --- | --- | --- | --- | --- | --- |
| | RE10K | LLFF | DTU | CO3D | Mip360 |
| ZipNeRF [47] | 20.77 | 17.23 | 9.18 | 14.34 | 12.77 |
| ZeroNVS [13] | 19.11 | 15.91 | 16.71 | 17.13 | 14.44 |
| ReconFusion [7] | 25.84 | 21.34 | 20.74 | 19.59 | 15.50 |
| CAT3D [8] | 26.78 | 21.58 | 22.02 | 20.57 | 16.62 |
| SEVA | 27.95 | 21.88 | 22.68 | 21.88 | 17.82 |

Table 4. **PSNR↑ on 3DGS renderings for set NVS.** Results are reported on the ReconFusion [7] split with $P = 3$.

tions, SEVA obtains a clear performance lead. On OO3D and GSO with $P = 1$, although the performances of SEVA and previous state-of-the-art method [19] are similar, we qualitatively observe more photorealistic and sharper output from our model.

**Qualitative comparison.** Fig. 6 top panel shows a qualitative comparison with diverse baselines. For small-viewpoint set NVS, the output from SEVA with the best scale exhibits desirable alignment with the ground truth while being more photorealistic in details. Compared with LVSM [11] on the P split of RE10K, SEVA produces sharper images, also corroborating that lower PSNR arises from scale ambiguity rather than interpolation quality. Similar trends hold when compared to Long-LRM [18] on DL3DV with $P = 32$. For large-viewpoint set NVS, we compare with DepthSplat [45] on DL3DV with $P = 3$. Depth-Splat fails to produce reasonable results when the viewpoint change is too large and falls short in overall visual quality.

**Comparison of 3D reconstruction.** To enable a direct quantitative comparison with prior works [7, 8], we adopt the few-view 3D reconstruction pipeline described in [8]. For each scene, we first generate 8 videos conditioned on the same input views following different camera paths,

summing into 720 generated views. Then, both the input views and generated views are distilled into a 3DGS-MCMC [50] representation without point cloud initialization. We optimize the camera parameters and apply LPIPS loss [51] during the distillation. Finally, we render the distilled 3D model on the test views and report the performance in Tab. 4. SEVA shows a consistent performance lead.

### 4.3. Trajectory NVS

In this section, we focus on qualitative demonstration, given that trajectory NVS is an underexplored task. We then compare against prior arts both qualitatively and quantitatively.

**Qualitative results.** Fig. 2 presents qualitative results, illustrating trajectories of varying complexities with different numbers of input views across diverse types, including object-centric scene-level, scene-level, real-world, and text-prompted from image diffusion models [25], *etc.*

In the single-view regime (*i.e.*, $P = 1$), we manually craft a set of common camera movements/effects, *e.g.*, look-at 360, spiral, panning, zoom-in, zoom-out, dolly zoom, *etc.* We observe that SEVA generalizes to a wide range of images and demonstrates accurate camera-following capacity. Excitingly, our model derives reasonable output with a dolly zoom effect (the second row of Fig. 2). In the FERN
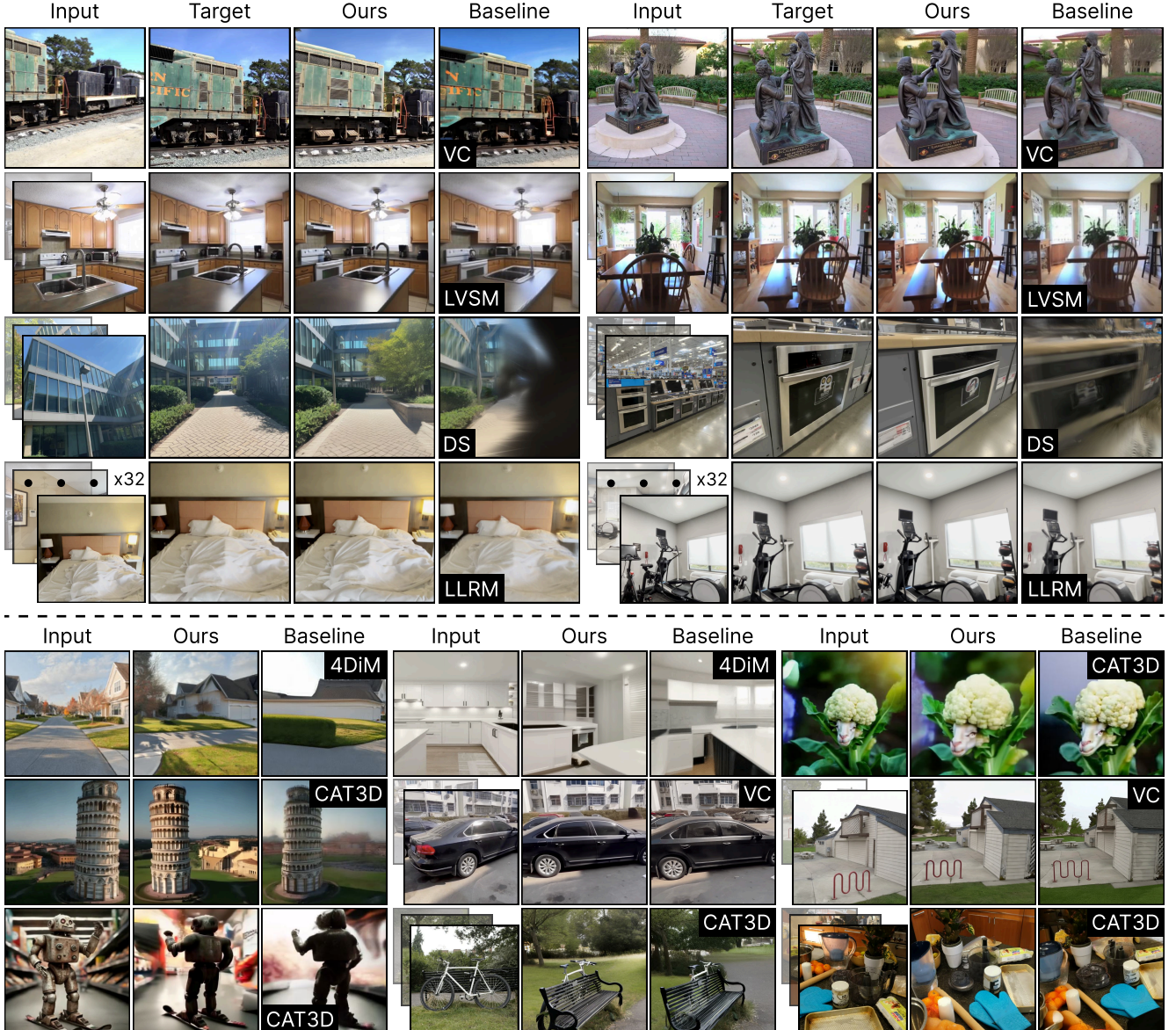
Figure 6. **SOTA comparison** on set NVS (top) and trajectory NVS (bottom) across varying numbers of input views. We compare with open-source approaches—ViewCrafter [9] (VC) and DepthSplat [45] (DS)—as well as proprietary ones including LVSM [11], Long-LRM [18] (LLRM), 4DiM [12], and CAT3D [8]. When the input comprises multiple views, we arrange them so that the view closest to the target is placed at the top of each set.

scene from the third row of Fig. 2, our model demonstrates its ability to generate plausible outputs even when moving close to or passing through an object—despite never being explicitly trained for such scenarios. This highlights the expressiveness of our model. An extensive sweeping of camera movements on 4 types of images is provided in Figs. 15 to 18.

In the sparse-view regime with few input views (*i.e.*, $1 < P \leqslant 8$), we observe that SEVA demonstrates strong generalization to in-the-wild real-world images and versatility in adapting to different numbers of input views. The

output forms a smooth trajectory video with subtle temporal flickering, revealing its capacity to interpolate between views smoothly. In the last row of Fig. 2, our model generates plausible results at the end of the trajectory—an area unseen in the input observations—demonstrating its strong generation capacity. In the semi-dense-view regime (*i.e.*, $P > 9$), we similarly find that SEVA is surprisingly able to produce a smooth trajectory video with minimal artifacts. Please check the website for video results.

**Qualitative comparison.** Fig. 6 bottom panel presents a qualitative comparison with diverse baselines. In the single-

| Method | split | small-viewpoint | | | large-viewpoint | | | | |
| | | V [9] | | | O | | | | |
| | dataset | RE | CO3D | T&T | RE | DTU | WR | DL | T&T |
|---|---|---|---|---|---|---|---|---|---|
| MotionCtrl [46] | | 16.29 | 15.46 | 13.29 | - | - | - | - | - |
| DepthSplat [45] | | 19.24 | 13.23 | 14.35 | 25.23 | 14.68 | 12.45 | 11.32 | 9.11 |
| ViewCrafter [9] | | 22.04 | 18.96 | 18.07 | 26.54 | 18.99 | 13.44 | 11.45 | 9.68 |
| SEVA | | 18.56 | 18.40 | 15.16 | 27.34 | 19.99 | 17.79 | 15.76 | 11.92 |
| SEVA (+ temp.) | | 18.62 | 18.43 | 15.13 | 27.36 | 20.19 | 17.93 | 15.78 | 11.99 |

Table 5. **PSNR↑ on trajectory NVS.** *temp.* denotes optional temporal pathway. RE, WR, and DL denotes RE10K, WRGBD, and DL3DV, respectively. For the V [9] split, $P = 1$ with unit length swept; for the O split, $P = 3$. Underlined numbers are run by us using the officially released code.

| Method | samples | | 3DGS | video |
| | PSNR↑ | TSED↓ | PSNR↑ | MS↑ |
|---|---|---|---|---|
| SEVA (one-pass) | 15.73 | 115.1 | 16.03 | 95.39 |
| SEVA (two-pass: nearest) | 13.74 | 120.9 | 14.21 | 94.71 |
| SEVA (two-pass: gt + nearest) | 15.58 | 116.2 | 15.96 | 95.22 |
| SEVA (two-pass: gt + interp) | 15.66 | 120.1 | 15.98 | 95.56 |
| SEVA | 15.76 | 116.7 | 16.11 | 95.76 |
| SEVA (+ temp.) | 15.78 | 109.0 | 16.17 | 95.77 |

Table 6. **3D consistency (TSED↓ and PSNR↑) and temporal quality (MS↑) on trajectory NVS.** SEVA uses interp procedural sampling by default. *temp.* denotes the optional temporal pathway. *MS* denotes motion smoothness from VBench [52]. Results are reported on our split of DL3DV with $P = 3$.

view regime (*i.e.*, $P = 1$), we compare to 4DiM [12] and CAT3D [8]. We observe more photo-realistic and sharper output from our model, especially in the background area for object-centric scenes. 4DiM outputs tend to be cartoonish and over-simplistic, given that the model is only trained on RE10K. In the sparse-view regime with few input views (*i.e.*, $P = 3$), we compare with CAT3D and observe that our model demonstrates more photo-realistic textures, especially in the background. For start-end-view interpolation considered in ViewCrafter [9] with $P = 2$, our model produces smooth transitions across trajectories, although it exhibits slight flickering between adjacent frames, particularly in regions with significant high-frequency detail.

**Quantitative comparison.** We use the same input views as in the set NVS for each split. We use all frames from each scene as target views such that they form a smoothly transitioning trajectory video, *i.e.*, $\mathbf{I}^{\text{tgt}} \sim \mathcal{V}$. We use PSNR as metrics and compare with baselines in Tab. 5.

For small-viewpoint trajectory NVS, Tab. 5 compares SEVA with baselines on PSNR. SEVA performs favorably against other methods in V split with $P = 1$. The performance lead of ViewCrafter is mainly attributed to its training on high-resolution images. For large-viewpoint trajectory NVS with $P = 3$, SEVA consistency sets new state-
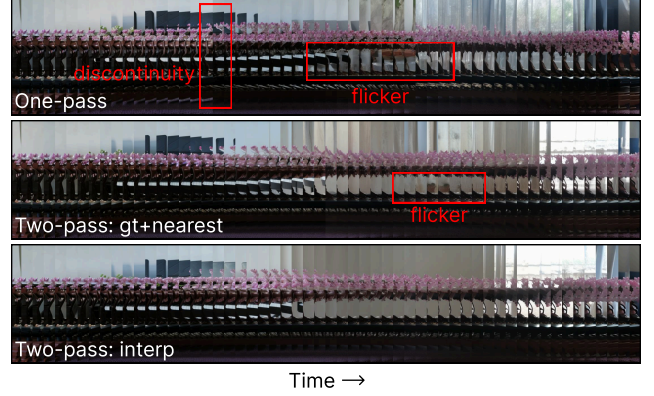


Figure 7. **Temporal quality.** Vertical slices of a rendered novel camera path on the BONSAI scene from Mip-NeRF360 [42] illustrate the temporal quality across adjacent viewpoints. One-pass or gt + nearest procedural sampling results in notable flickering, whereas interp procedural sampling ensures temporally smooth rendering.

of-the-art results. Applying the temporal pathway further boosts performance and improves smoothness, indicating the benefits of the gated architecture.

**Ablation on two-pass procedural sampling.** We conduct an ablation study comparing the default interp procedural sampling with one-pass sampling and alternative procedural sampling strategies.

Quantitatively, beyond evaluating PSNR on individual views, we assess 3D consistency using the PSNR on 3D renderings of that same view and SED [12, 53] score. To compute the SED score, we first apply SIFT [54] to detect keypoints in two images. For each keypoint in the first image, we determine its corresponding epipolar line in the second image and measure the shortest distance to its match. Additionally, we report Motion Smoothness (MS) from VBench [52], a benchmark designed to evaluate temporal coherence in video generative models. As shown in Tab. 6, interp procedural sampling demonstrates a clear advantage over its alternatives, with the integration of the temporal pathway further reinforcing its superiority.

Qualitative comparisons in Fig. 7 show that one-pass sampling introduces visible temporal flickering and abrupt visual changes. In contrast, interp produces the smoothest transitions, outperforming gt + nearest and mitigating noticeable flickering.

### 4.4. Long-Trajectory NVS

Fig. 8 presents a qualitative demonstration of NVS over a long trajectory of up to 1000 frames. As the camera orbits the TELEPHONE BOOTH for multiple rounds, the generated views in each round from similar viewpoints can be drastically different since they are far away from each other temporally. With the memory bank maintaining previously generated anchors, SEVA achieves robust 3D con-
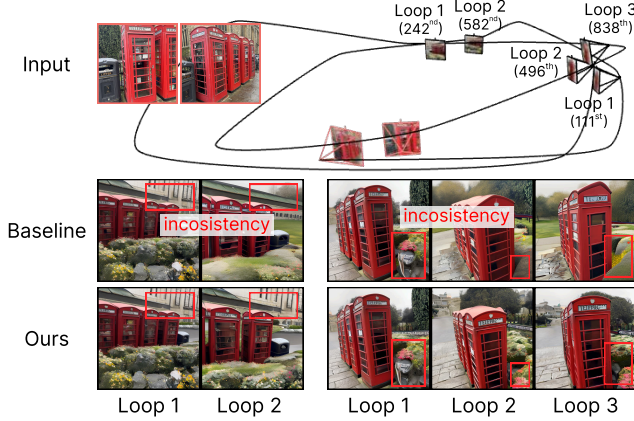
Figure 8. **Long-range 3D consistency.** We visualize samples following a camera path looping three times around the TELEPHONE-BOOTH scene. Lookup using spatial neighbors from the memory bank (ours) notably improves view consistency and reduces artifacts in recurring locations across different loops, compared to lookup using temporal neighbors (baseline).



Figure 9. **Generation quality on the number of input views.** PSNR↑ (top) and Image Quality↑ (bottom) on set NVS. Results are reported on our split of T&T. Extending $T$ to more input views in a zero-shot manner produces more consistent samples in the semi-dense-view regime. Dense 3DGS denotes results of [3] with full views.

sistency for long-trajectory NVS, *e.g.*, the building in front of and the plantation after the booth. Comparing it to using temporal nearest anchors previously generated, using spatially nearest ones demonstrates a clear advantage. The memory mechanism has been concurrently explored in previous works [9, 55], leveraging explicit intermediate 3D representations such as dense point clouds predicted by DUSt3R [56]. In contrast, our model demonstrates greater robustness and generalizability to in-the-wild data, as it is not constrained by the quality of DUSt3R's output, which often becomes unreliable in quality for data outside of its training domain, *e.g.*, text-prompted images.

## 4.5. Discussions

**Zero-shot generalization of context window length $T$.** We surprisingly find our model, though only trained on $T = 21$ frames, can generalize reasonably to larger $T$ during sampling *in the semi-dense-view regime*. On our split of T&T for set NVS, we evaluate the predictions against ground truth in both sparse-view (*i.e.*, $1 \leq P \leq 8$) and semi-dense-view regime (*i.e.*, $9 \leq P$) using PSNR↑ and Image Quality↑ [52]. Image Quality refers to the distortion (*e.g.*, over-exposure, noise, blur) presented in the generated image. We experiment with different sampling strategies: one-pass sampling zero-shot extending the context window length $T$; two-pass procedural sampling by first generating anchor views using nearest−$K$ ($K < T$) input views and then interpolating anchor views into target views.

Our results are shown in Fig. 9. Procedural sampling with the nearest−$K$ anchor views plateau after taking $K$ views as input, indicating inefficiencies in procedural sampling and an inability to effectively utilize all available input views when $P > T$. Conversely, the metrics steadily im-

prove with respect to the number of input frames for one-pass sampling with $T$ extending to $P + Q$ in a zero-shot manner. However, we observe that this generalization fails in the sparse-view regime, resulting in blurry samples, as indicated by the low Image Quality when $P < 9$ and qualitative samples when $P = 3$. In the semi-dense-view setting, although quantitative metrics show minimal differences between one-pass and procedural sampling, we consistently observe that one-pass produces more 3D-consistent samples, as illustrated in the bottom-right figure.

**Zero-shot generalization of image resolution.** Surprisingly, we find our model, despite being trained only on square images with $H = W = 576$, generalizes well to different image resolution during sampling, similar to [57]. As shown in Fig. 10, SEVA can produce high-quality results in both portrait ($16 : 9$) and landscape ($9 : 16$) orientations of different image resolutions.

**Guidance scale on generation uncertainty.** We employ classifier-free guidance [58] (CFG) to enhance sampling quality. Empirically, we find that the CFG scale, a hyper-parameter at test time, has a significant impact on the final result [22], as shown in Fig. 11. Specifically, the optimal CFG scale is strongly correlated with the inherent uncertainty of the generation. When uncertainty is high (top row), a higher CFG scale (*e.g.*, 5) is preferable to prevent excessive blurriness in the generated samples. Conversely, when uncertainty is low (bottom row), a lower CFG scale (*e.g.*, 3) helps avoid oversaturation. In practice, setting the CFG scale between 2 and 5 consistently produces high-quality results across all our samples.
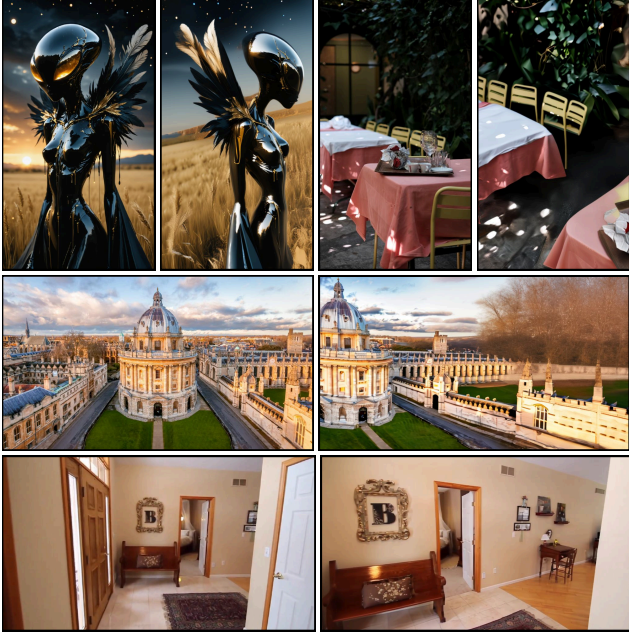
Figure 10. **Generation quality on different image resolutions.** Our model generalizes to different image resolution of varying aspect ratios, including both portrait (top) and landscape orientations (bottom). Results are presented as a pair of the input view and the target views.
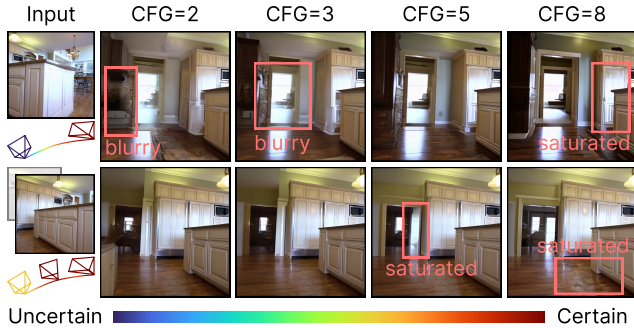


Uncertain ▰▰▰▰▰▰▰▰▰▰▰▰ Certain

Figure 11. **Generation uncertainty on CFG.** The CFG scale should be increased as generation uncertainty rises. For single-view conditioning (top), a higher CFG scale is typically required, whereas few-view conditioning (bottom) benefits from a lower scale.

**Sampling diversity of unseen areas.** Fig. 12 demonstrates the capability of the model to generate diverse and plausible predictions for unseen regions of input observations. In the first row, the input view depicts a frontal view of a classical statue. We sample multiple back views by varying the random seeds, producing distinct yet coherent interpretations of the unseen geometry and texture while preserving fidelity to the input. Similarly, in the second row, the model generates multiple plausible continuations of the scene given an input view of a scenic road, each reflecting unique variations in environmental and structural



Figure 12. **Generation diversity in unseen regions.** Our model generates diverse samples by varying randomization seeds during the sampling process.

details. These results highlight the model's ability to synthesize realistic and diverse outputs for occluded or ambiguous regions.

## 5. Conclusion

We present STABLE VIRTUAL CAMERA (SEVA), a generalist diffusion model for novel view synthesis that balances large viewpoint changes and smooth interpolation while supporting flexible input and target configurations. By designing a diffusion-based architecture without 3D representation, a structured training strategy, and a two-pass procedural sampling approach, SEVA achieves 3D consistent rendering across diverse NVS tasks. Extensive benchmarking demonstrates its superiority over existing methods, with strong generalization to real-world scenes. For broader impact and limitations, please refer to the appendix.

# References

[1] 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. In *European Conference on Computer Vision*, pages 405–421, 2020. 2

[2] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Trans. Graph.*, 2022.

[3] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3D gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4), July 2023. 2, 4, 11, 6

[4] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelNeRF: Neural radiance fields from one or few images. In *CVPR*, 2021. 2

[5] Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. Dreamfusion: Text-to-3D using 2D diffusion. *arXiv*, 2022. 2

[6] 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. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023. 2

[7] Rundi Wu, Ben Mildenhall, Philipp Henzler, Keunhong Park, Ruiqi Gao, Daniel Watson, Pratul P Srinivasan, Dor Verbin, Jonathan T Barron, Ben Poole, et al. Reconfusion: 3d reconstruction with diffusion priors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21551–21561, 2024. 2, 4, 7, 8, 3, 5

[8] Ruiqi Gao*, Aleksander Holynski*, Philipp Henzler, Arthur Brussee, Ricardo Martin-Brualla, Pratul P. Srinivasan, Jonathan T. Barron, and Ben Poole*. Cat3d: Create anything in 3d with multi-view diffusion models. *arXiv*, 2024. 2, 4, 5, 6, 7, 8, 9, 10

[9] Wangbo Yu, Jinbo Xing, Li Yuan, Wenbo Hu, Xiaoyu Li, Zhipeng Huang, Xiangjun Gao, Tien-Tsin Wong, Ying Shan, and Yonghong Tian. ViewCrafter: Taming video diffusion models for high-fidelity novel view synthesis. *arXiv preprint arXiv:2409.02048*, 2024. 2, 4, 7, 8, 9, 10, 11, 3, 5

[10] Zhouxia Wang, Ziyang Yuan, Xintao Wang, Yaowei Li, Tianshui Chen, Menghan Xia, Ping Luo, and Ying Shan. Motionctrl: A unified and flexible motion controller for video generation. In *ACM SIGGRAPH 2024 Conference Papers*, pages 1–11, 2024. 2, 4

[11] Haian Jin, Hanwen Jiang, Hao Tan, Kai Zhang, Sai Bi, Tianyuan Zhang, Fujun Luan, Noah Snavely, and Zexiang Xu. Lvsm: A large view synthesis model with minimal 3d inductive bias, 2024. 2, 7, 8, 9, 4

[12] Daniel Watson, Saurabh Saxena, Lala Li, Andrea Tagliasacchi, and David J Fleet. Controlling space and time with diffusion models. *arXiv preprint arXiv:2407.07860*, 2024. 2, 7, 8, 9, 10, 3, 4

[13] Kyle Sargent, Zizhang Li, Tanmay Shah, Charles Herrmann, Hong-Xing Yu, Yunzhi Zhang, Eric Ryan Chan, Dmitry Lagun, Li Fei-Fei, Deqing Sun, et al. Zeronvs: Zero-shot

[14] Virtual camera system. https://en.wikipedia.org/wiki/Virtual_camera_system. Accessed: 5-Mar-2025. 2

[15] How weta digital's virtual camera transforms cinematography. https://www.youtube.com/watch?v=I-ctZRxLEN4. Accessed: 5-Mar-2025.

[16] David Charatan, Sizhe Li, Andrea Tagliasacchi, and Vincent Sitzmann. pixelsplat: 3d gaussian splats from image pairs for scalable generalizable 3d reconstruction. In *arXiv*, 2023. 2, 7, 8, 3, 4

[17] Yuedong Chen, Haofei Xu, Chuanxia Zheng, Bohan Zhuang, Marc Pollefeys, Andreas Geiger, Tat-Jen Cham, and Jianfei Cai. Mvsplat: Efficient 3d gaussian splatting from sparse multi-view images. In *European Conference on Computer Vision*, pages 370–386. Springer, 2025. 2, 7, 8, 4

[18] Chen Ziwen, Hao Tan, Kai Zhang, Sai Bi, Fujun Luan, Yicong Hong, Li Fuxin, and Zexiang Xu. Long-lrm: Long-sequence large reconstruction model for wide-coverage gaussian splats. *arXiv preprint arXiv:2410.12781*, 2024. 2, 7, 8, 9, 3, 4

[19] Vikram Voleti, Chun-Han Yao, Mark Boss, Adam Letts, David Pankratz, Dmitrii Tochilkin, Christian Laforte, Robin Rombach, and Varun Jampani. SV3D: Novel multi-view synthesis and 3D generation from a single image using latent video diffusion. In *European Conference on Computer Vision*, 2024. 2, 4, 5, 7, 8, 3

[20] 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. *Communications of the ACM*, 2021. 4

[21] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 4, 2

[22] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023. 4, 5, 11

[23] Hao He, Yinghao Xu, Yuwei Guo, Gordon Wetzstein, Bo Dai, Hongsheng Li, and Ceyuan Yang. Cameractrl: Enabling camera control for text-to-video generation. *arXiv preprint arXiv:2404.02101*, 2024.

[24] Yingqing He, Tianyu Yang, Yong Zhang, Ying Shan, and Qifeng Chen. Latent video diffusion models for high-fidelity long video generation. *arXiv*, 2022. 4, 6

[25] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10684–10695, June 2022. 4, 8

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

[27] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceed-*

13

ings of the IEEE conference on computer vision and pattern recognition, pages 770–778, 2016. 4

[28] P Goyal. Accurate, large minibatch sg d: training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017. 4

[29] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22563–22575, 2023. 5

[30] Julius Plucker. Xvii. on a new geometry of space. *Philosophical Transactions of the Royal Society of London*, (155):725–791, 1865. 5

[31] Chuanxia Zheng and Andrea Vedaldi. Free3D: Consistent novel view synthesis without 3D representation. *arXiv preprint arXiv:2312.04551*, 2023. 5

[32] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 5, 3

[33] Emiel Hoogeboom, Jonathan Heek, and Tim Salimans. simple diffusion: End-to-end diffusion for high resolution images. *arXiv preprint arXiv:2301.11093*, 2023. 5

[34] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *ICLR*, 2024. 5

[35] Tong Wu, Jiarui Zhang, Xiao Fu, Yuxin Wang, Jiawei Ren, Liang Pan, Wayne Wu, Lei Yang, Jiaqi Wang, Chen Qian, et al. Omniobject3D: Large-vocabulary 3D object dataset for realistic perception, reconstruction and generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 803–814, 2023. 7, 3

[36] Laura Downs, Anthony Francis, Nate Koenig, Brandon Kinman, Ryan Hickman, Krista Reymann, Thomas B McHugh, and Vincent Vanhoucke. Google Scanned Objects: A high-quality dataset of 3D scanned household items. In *2022 International Conference on Robotics and Automation (ICRA)*, pages 2553–2560. IEEE, 2022. 7, 3

[37] Ben Mildenhall, Pratul P Srinivasan, Rodrigo Ortiz-Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi, Ren Ng, and Abhishek Kar. Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Transactions on Graphics (ToG)*, 38(4):1–14, 2019. 7, 3

[38] Rasmus Jensen, Anders Dahl, George Vogiatzis, Engin Tola, and Henrik Aanæs. Large scale multi-view stereopsis evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 406–413, 2014. 7, 3

[39] Jeremy Reizenstein, Roman Shapovalov, Philipp Henzler, Luca Sbordone, Patrick Labatut, and David Novotny. Common objects in 3d: Large-scale learning and evaluation of real-life 3d category reconstruction. In *International Conference on Computer Vision*, 2021. 7, 3

[40] Hongchi Xia, Yang Fu, Sifei Liu, and Xiaolong Wang. Rgbd objects in the wild: Scaling real-world 3d object learning from rgb-d videos, 2024. 7, 3, 5

[41] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: Learning view synthesis using multiplane images. *arXiv preprint arXiv:1805.09817*, 2018. 7, 3

[42] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-NeRF: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5855–5864, 2021. 7, 10, 3

[43] Lu Ling, Yichen Sheng, Zhi Tu, Wentian Zhao, Cheng Xin, Kun Wan, Lantao Yu, Qianyu Guo, Zixun Yu, Yawen Lu, et al. Dl3dv-10k: A large-scale scene dataset for deep learning-based 3d vision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22160–22169, 2024. 7, 3, 5

[44] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Transactions on Graphics (ToG)*, 36(4):1–13, 2017. 7, 3, 5

[45] Haofei Xu, Songyou Peng, Fangjinhua Wang, Hermann Blum, Daniel Barath, Andreas Geiger, and Marc Pollefeys. Depthsplat: Connecting gaussian splatting and depth. *arXiv preprint arXiv:2410.13862*, 2024. 7, 8, 9, 10, 4, 5

[46] Tsai-Shien Chen, Chieh Hubert Lin, Hung-Yu Tseng, Tsung-Yi Lin, and Ming-Hsuan Yang. Motion-conditioned diffusion model for controllable video synthesis. *arXiv preprint arXiv:2304.14404*, 2023. 7, 8, 10, 4, 5

[47] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Zip-nerf: Anti-aliased grid-based neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19697–19705, 2023. 8, 5

[48] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018. 7, 6

[49] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. 7, 6

[50] Shakiba Kheradmand, Daniel Rebain, Gopal Sharma, Weiwei Sun, Jeff Tseng, Hossam Isack, Abhishek Kar, Andrea Tagliasacchi, and Kwang Moo Yi. 3d gaussian splatting as markov chain monte carlo. *arXiv preprint arXiv:2404.09591*, 2024. 8

[51] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric, 2018. 8

[52] Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, Yaohui Wang, Xinyuan Chen, Limin Wang, Dahua Lin, Yu Qiao, and Ziwei Liu. VBench: Comprehensive benchmark suite for video generative models. In

*Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. 10, 11

[53] Jason J Yu, Fereshteh Forghani, Konstantinos G Derpanis, and Marcus A Brubaker. Long-term photometric consistent novel view synthesis with diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7094–7104, 2023. 10

[54] Huiyu Zhou, Yuan Yuan, and Chunmei Shi. Object tracking using sift features and mean shift. *Computer vision and image understanding*, 113(3):345–352, 2009. 10

[55] Baorui Ma, Huachen Gao, Haoge Deng, Zhengxiong Luo, Tiejun Huang, Lulu Tang, and Xinlong Wang. You see it, you got it: Learning 3d creation on pose-free videos at scale. 2024. 11

[56] Shuzhe Wang, Vincent Leroy, Yohann Cabon, Boris Chidlovskii, and Jerome Revaud. Dust3r: Geometric 3d vision made easy. In *CVPR*, 2024. 11

[57] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 11

[58] Jonathan Ho and Tim Salimans. Classifier-Free Diffusion Guidance. *arXiv:2207.12598*, 2022. 11

[59] 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. 2

[60] Yang Song and Stefano Ermon. Generative Modeling by Estimating Gradients of the Data Distribution. *arXiv:1907.05600*, 2019. 2

[61] Mohammed Suhail, Carlos Esteves, Leonid Sigal, and Ameesh Makadia. Light field neural rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8269–8279, 2022. 2

[62] Stanislaw Szymanowicz, Chrisitian Rupprecht, and Andrea Vedaldi. Splatter image: Ultra-fast single-view 3d reconstruction. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10208–10217, 2024. 2

[63] Yicong Hong, Kai Zhang, Jiuxiang Gu, Sai Bi, Yang Zhou, Difan Liu, Feng Liu, Kalyan Sunkavalli, Trung Bui, and Hao Tan. LRM: Large reconstruction model for single image to 3D. *arXiv preprint arXiv:2311.04400*, 2023. 2

[64] Daniel Watson, William Chan, Ricardo Martin-Brualla, Jonathan Ho, Andrea Tagliasacchi, and Mohammad Norouzi. Novel view synthesis with diffusion models, 2022. 2

[65] 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. *International Conference on Computer Vision*, 2023.

[66] Ruoxi Shi, Hansheng Chen, Zhuoyang Zhang, Minghua Liu, Chao Xu, Xinyue Wei, Linghao Chen, Chong Zeng, and Hao Su. Zero123++: a single image to consistent multi-view diffusion base model. *arXiv preprint arXiv:2310.15110*, 2023. 2

[67] Yuan Liu, Cheng Lin, Zijiao Zeng, Xiaoxiao Long, Lingjie Liu, Taku Komura, and Wenping Wang. SyncDreamer: Generating multiview-consistent images from a single-view image. *arXiv preprint arXiv:2309.03453*, 2023. 2

[68] Antoine Mercier, Ramin Nakhli, Mahesh Reddy, Rajeev Yasarla, Hong Cai, Fatih Porikli, and Guillaume Berger. HexaGen3D: Stablediffusion is just one step away from fast and diverse Text-to-3D generation. *arXiv preprint arXiv:2401.07727*, 2024. 2

[69] Sherwin Bahmani, Ivan Skorokhodov, Victor Rong, Gordon Wetzstein, Leonidas Guibas, Peter Wonka, Sergey Tulyakov, Jeong Joon Park, Andrea Tagliasacchi, and David B Lindell. 4d-fy: Text-to-4d generation using hybrid score distillation sampling. *arXiv preprint arXiv:2311.17984*, 2023. 2

[70] Petar Veličković, Christos Perivolaropoulos, Federico Barbero, and Razvan Pascanu. softmax is not enough (for sharp out-of-distribution). *arXiv preprint arXiv:2410.01104*, 2024. 6

[71] Kiwhan Song, Boyuan Chen, Max Simchowitz, Yilun Du, Russ Tedrake, and Vincent Sitzmann. History-guided video diffusion, 2025. 6

[72] Boyuan Chen, Diego Marti Monso, Yilun Du, Max Simchowitz, Russ Tedrake, and Vincent Sitzmann. Diffusion forcing: Next-token prediction meets full-sequence diffusion, 2024. 6

[73] Stability AI. Stable diffusion 3.5. https://stability.ai/news/introducing-stable-diffusion-3-5, 2024. 7, 8

# Contents

## A. Broader Impact and Limitations

**Broader Impact.** SEVA significantly advances immersive 3D experiences by synthesizing realistic and temporally consistent views from sparse camera inputs, addressing key limitations in NVS. Inspired by James Cameron's pioneering Virtual Camera technology—which enabled filmmakers to intuitively navigate virtual environments and visualize precise camera trajectories—our generative AI-driven model similarly allows users to create intricate, controllable camera paths without the typical complexity of dense view captures or explicit 3D reconstructions. By generalizing across arbitrary viewpoint changes and enabling temporally smooth rendering without NeRF distillation, our approach simplifies the NVS pipeline, enhancing accessibility for content creators, developers, and researchers. This facilitates applications ranging from virtual cinematography and gaming to digital heritage preservation, substantially broadening the usability and scalability of NVS.

**Limitations.** The performance of SEVA is constrained by the scope of its training data, resulting in reduced quality for certain types of scenes. Specifically, input images featuring humans, animals, or dynamic textures (*e.g.*, water surfaces) typically lead to degraded outputs. Additionally, highly ambiguous scenes or complex camera trajectories pose challenges; for instance, trajectories that intersect with objects or surfaces may cause noticeable flickering artifacts. Similar issues arise for extremely irregularly shaped objects or when target viewpoints significantly diverge from the provided input viewpoints.

## B. Related Work

**Novel view synthesis.** While traditional NVS has been studied for nearly several decades, it has recently achieved remarkable success with the help of techniques such as NeRF [1, 59] and diffusion models [21, 60]. Using these techniques, there are broadly two ways of generating novel views : 1) estimate a 3D representation using multiple sparse input views, then regress the novel views from this intermediate representation, 2) directly estimate the novel views from the sparse input views, either in a single shot in a feed-forward manner, or in multiple sampling steps using diffusion models.

**Feed-forward models.** Approaches like LFNR [61] and LVSM [11] directly generate target views and leverage data-driven learning to capture 3D inductive biases. While often efficient, these methods struggle with the inherent diversity of generative NVS, limiting their capacity to model multiple plausible solutions. In contrast, our approach frames generative NVS through a diffusion perspective, enabling us to sample diverse, plausible solutions during inference,

thereby addressing ambiguities and enhancing generation capacity.

**Intermediate representation models.** Techniques such as NeRF [1] and Gaussian Splatting [3] have made significant progress on per-scene optimization from input views by building 3D representations efficiently. Several works show that these representations can then be used to regress novel views. pixelNeRF [4] builds a NeRF from multiple input views; Splatter Image [62], pixelSplat [16], and MVSplat [17] build a 3D representation using Gaussian Splatting; LRM [63] builds a triplane representation. However, these optimization-based methods cannot creatively synthesize missing regions, and rely on tens, if not hundreds, of posed input images which limits their practicality in real-world applications.

**Diffusion-based models.** Our work falls within this category, where target novel views are generated in multiple steps through a denoising diffusion process [21, 60]. As mentioned earlier, existing diffusion-based methods can be divided into two main types: *image models* and *video models*.

Image models are designed to synthesize distant viewpoints [64–66]. However, these early practices only generate one viewpoint at a time, and lack multi-view consistency, often resulting in jittery and inconsistent samples when generating along a camera trajectory. Works such as MVDream [26], SyncDreamer [67] and HexGen3D [68] generate multiple fixed views simultaneously. However, these models only generate specific views given a conditional image, not arbitrary viewpoints.

To obtain consistent 3D objects, these models necessitate NeRF distillation, either through Score Distillation Samling (SDS) [5, 13] or directly upon completely sampled images [7, 8].

Video models can produce smooth video sequences by maintaining certain constraints relative to the input views [19]. However, they are generally limited to smaller camera motions due to the natural frame rate in video training. Some works use video diffusion models to generate 4D scenes [69]. But in those works, the video diffusion models do not contribute to the consistency of the 3D object itself, that part is handled by image-based diffusion models such as MVDream.

## C. Benchmark

We collect 10 commonly used datasets to benchmark NVS, encompassing a diverse range of scene distributions and complexities, shown in Tab. 7.

| | type | split | #scene | $(I^{inp}, I^{tgt}) \sim \mathcal{V}$ | $P$ | $\mathcal{D}_{CLIP}(I)$ |
|---|---|---|---|---|---|---|
| **Small-viewpoint NVS** | | | | | | |
| OmniObject3D [35] | ▣ | O (dynamic orbit) | 308 | ✓ | 3 | 0.11 |
| GSO [36] | ▣ | O (dynamic orbit) | 300 | ✓ | 3 | 0.11 |
| RealEstate10K [41] | ▲ | D [12] | 128 | ✓ | 1 | 0.09 |
| | | R [7] | 10 | ✓ | 1 | 0.08 |
| | | | | | 3 | 0.03 |
| | | P [16] | 6474 | ✓ | 2 | 0.04 |
| | | V [9] | 10 | ✓ | 2 | 0.11 |
| LLFF [37] | ▲ | R [7] | 8 | ✓ | 1 | 0.04 |
| | | | | | 3 | 0.03 |
| DTU [38] | ▲ | R [7] | 15 | ✓ | 1 | 0.07 |
| | | | | | 3 | 0.06 |
| CO3D [39] | ▲ | R [7] | 20 | ✓ | 3 | 0.09 |
| | | V [9] | 10 | ✓ | 2 | 0.09 |
| WildRGB-D [40] | ▲ | $O_e$ (1/3 orbit) | 20 | ✓ | 3 | 0.07 |
| | | $O_h$ (full orbit) | | | 6 | 0.11 |
| Mip-NeRF360 [42] | ▲ | R [7] | 9 | ✗ | 6 | 0.11 |
| DL3DV-140 [43] | ▲ | O | 10 | ✓ | 6 | 0.10 |
| | | L [18] | 140 | ✓ | 32 | 0.05 |
| Tanks and Temples [44] | ▲ | V [9] | 22 | ✓ | 2 | 0.10 |
| | | L [18] | 2 | ✓ | 32 | 0.10 |
| **Large-viewpoint NVS** | | | | | | |
| OmniObject3D [35] | ▣ | S [19] (dynamic orbit) | 308 | ✓ | 1 | 0.16 |
| GSO [36] | ▣ | S [19] (dynamic orbit) | 300 | ✓ | 1 | 0.18 |
| CO3D [39] | ▲ | R [7] | 20 | ✓ | 1 | 0.15 |
| WildRGB-D [40] | ▲ | $O_h$ (full orbit) | 20 | ✓ | 1 | 0.19 |
| | | | | | 3 | 0.14 |
| Mip-NeRF360 [42] | ▲ | R [7] | 9 | ✗ | 1 | 0.19 |
| | | | | | 3 | 0.13 |
| DL3DV-140 [43] | ▲ | O | 10 | ✓ | 1 | 0.21 |
| | | | | | 3 | 0.12 |
| Tanks and Temples [44] | ▲ | O | 2 | ✓ | 1 | 0.21 |
| | | | | | 3 | 0.18 |
| | | | | | 6 | 0.16 |
| | | | | | 9 | 0.14 |

Table 7. **Statistics for NVS benchmark**. We consider 10 publicly available datasets commonly used for evaluating NVS, encompassing both object-level and scene-level data. Views from Mip-NeRF360 [42] derive from several disjoint captures following different camera trajectories, thus all views $(I^{inp}, I^{tgt}) \sim \mathcal{I}$. $P$ denotes the number of input views. Depending on the disparity between $I^{inp}$ and $I^{tgt}$, we group NVS tasks into small-viewpoint NVS (top panel) where target views are similar to input views and large-viewpoint NVS (bottom panel) where target views are more different to input views.

**Small-viewpoint *versus* large-viewpoint NVS.** In Sec. 4.1, we split NVS tasks into two categories: small-viewpoint and large-viewpoint NVS based on the disparity between $I^{inp}$ and $I^{tgt}$. Formally, for each target view, we consider the minimal distance between the CLIP [32] feature of that view and those of all input views. Averaging across all target views yields the CLIP distance, $\mathcal{D}_{CLIP}(I)$. Splits with $\mathcal{D}_{CLIP}(I) <= 0.11$ are grouped as small-view NVS, while those with $\mathcal{D}_{CLIP}(I) > 0.11$

are grouped as large-view NVS. We concrete in Tab. 7 a detailed task setup including the choice of datasets and splits (depending on which scenes from each dataset and which views from each scene are used).

**Choice of scenes.** We follow the choices of scenes for splits adopted from previous works. For our split, we use all scenes from the dataset without specification.

For the Tanks and Temples dataset, the 2 chosen scenes

**(a) LPIPS↓**

| Method | OO3D | GSO | RE10K | | | | | LLFF | | DTU | | CO3D | | WRGBD | | Mip360 | DL3DV | | T&T | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| split | O | O | D [12] | V [9] | P [16] | R [7] | | R [7] | | R [7] | | V [9] | R [7] | $O_e$ | $O_h$ | R [7] | O | L [18] | V [9] | L [18] |
| $P$ | 3 | 3 | 1 | 1 | 2 | 1 | 3 | 1 | 3 | 1 | 3 | 1 | 3 | 3 | 6 | 6 | 6 | 32 | 1 | 32 |
| **Regression-based models** | | | | | | | | | | | | | | | | | | | | |
| Long-LRM [18] | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 0.262 | - | 0.375 |
| MVSplat [17] | 0.411 | 0.387 | 0.224 | 0.237 | 0.128 | 0.254 | 0.142 | 0.542 | 0.497 | 0.386 | 0.310 | 0.634 | 0.614 | 0.504 | 0.643 | 0.556 | 0.527 | 0.425 | 0.519 | 0.568 |
| DepthSplat [45] | 0.404 | 0.372 | 0.217 | 0.245 | 0.119 | 0.236 | 0.177 | 0.530 | 0.465 | 0.369 | 0.304 | 0.618 | 0.603 | 0.499 | 0.530 | 0.534 | 0.481 | 0.404 | 0.462 | 0.528 |
| LVSM [11] | - | - | - | - | 0.098 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - |
| **Diffusion-based models** | | | | | | | | | | | | | | | | | | | | |
| MotionCtrl [46] | - | - | 0.500 | 0.386 | - | - | - | - | - | - | - | 0.443 | - | - | - | - | - | - | 0.473 | - |
| 4DiM [12] | - | - | 0.302 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - |
| ViewCrafter [9] | 0.427 | 0.379 | 0.220 | 0.178 | 0.203 | 0.287 | 0.164 | 0.620 | 0.435 | 0.485 | 0.272 | 0.324 | 0.513 | 0.324 | 0.639 | 0.464 | 0.558 | - | 0.283 | - |
| SEVA | 0.049 | 0.041 | 0.194 | 0.231 | 0.061 | 0.308 | 0.073 | 0.389 | 0.181 | 0.316 | 0.158 | 0.318 | 0.278 | 0.215 | 0.237 | 0.319 | 0.232 | 0.155 | 0.354 | 0.236 |

**(b) SSIM↑**

| Method | OO3D | GSO | RE10K | | | | | LLFF | | DTU | | CO3D | | WRGBD | | Mip360 | DL3DV | | T&T | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| split | O | O | D [12] | V [9] | P [16] | R [7] | | R [7] | | R [7] | | V [9] | R [7] | $O_e$ | $O_h$ | R [7] | O | L [18] | V [9] | L [18] |
| $P$ | 3 | 3 | 1 | 1 | 2 | 1 | 3 | 1 | 3 | 1 | 3 | 1 | 3 | 3 | 6 | 6 | 6 | 32 | 1 | 32 |
| **Regression-based models** | | | | | | | | | | | | | | | | | | | | |
| Long-LRM [18] | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 0.775 | - | 0.590 |
| MVSplat [17] | 0.554 | 0.621 | 0.788 | 0.769 | 0.869 | 0.812 | 0.857 | 0.283 | 0.358 | 0.576 | 0.624 | 0.403 | 0.370 | 0.405 | 0.368 | 0.312 | 0.487 | 0.512 | 0.394 | 0.314 |
| DepthSplat [45] | 0.636 | 0.689 | 0.801 | 0.745 | 0.887 | 0.820 | 0.824 | 0.299 | 0.396 | 0.601 | 0.638 | 0.429 | 0.402 | 0.436 | 0.417 | 0.324 | 0.513 | 0.564 | 0.413 | 0.359 |
| LVSM [11] | - | - | - | - | 0.906 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - |
| **Diffusion-based models** | | | | | | | | | | | | | | | | | | | | |
| MotionCtrl [46] | - | - | 0.267 | 0.587 | - | - | - | - | - | - | - | 0.502 | - | - | - | - | - | - | 0.384 | - |
| 4DiM [12] | - | - | 0.463 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - |
| ViewCrafter [9] | 0.538 | 0.647 | 0.792 | 0.798 | 0.710 | 0.806 | 0.830 | 0.146 | 0.454 | 0.542 | 0.671 | 0.641 | 0.483 | 0.465 | 0.376 | 0.354 | 0.469 | - | 0.563 | - |
| SEVA | 0.935 | 0.942 | 0.615 | 0.693 | 0.847 | 0.700 | 0.892 | 0.384 | 0.602 | 0.652 | 0.750 | 0.585 | 0.647 | 0.670 | 0.646 | 0.395 | 0.546 | 0.661 | 0.437 | 0.505 |

Table 8. **LPIPS↓ (top) and SSIM↑ (bottom) on small-viewpoint set NVS.** For all results with $P = 1$, we sweep the unit length for camera normalization due to the model's scale ambiguity. $O_e$ and $O_h$ denote the easy and hard split of our split. Underlined numbers are run by us using the officially released code.

in our (O) split are TRAIN and TRUCK. For the DL3DV-140 dataset, the 10 test scenes we choose in O split are:

- 165F5AF8BFE32F70595A1C9393A6E442ACF7AF019998275144F605B89A306557
- 341B4FF3DFD3D377D7167BD81F443BEDAFBFF003BF04881B99760FC0AEB69510
- 3BB3BB4D3E871D79EB71946CBAB1E3AFC7A8E33A661153033F32DEB3E23D2E52
- 3BB894D1933F3081134AD2D40E54DE5F0636BD8B502B0A8561873BB63B0DCE85
- 9E9A89AE6FED06D6E2F4749B4B0059F35CA97F848CEDC4A14345999E746F7884
- CD9C981EEB4A9091547AF19181B382698E9D9EEE0A838C7C9783A8A268AF6AEE
- D4FBEBA0168AF8FDDB2FC695881787AEDCD62F477C7DCEC9EBCA7B8594BBD95B
- E78F8CEBD2BD93D960BFAEAC18FAC0BB2524F15C44288903CD20B73E599E8A81
- ED16328235C610F15405FF08711EAF15D88A0503884F3A9CCB5A0EE69CB4ACB5

- F71AC346CD0FC4652A89AFB37044887EC3907D37D01D1CEB0AD28E1A780D8E03.

For the WildRGBD dataset, the 20 test scenes we choose in O split are: For the WildRGBD dataset, the 20 test scenes we select for the O split are:

- BALL/SCENE_563
- APPLE/SCENE_234
- MICROWAVE/SCENE_143
- SCISSOR/SCENE_489
- BUCKET/SCENE_294
- KEYBOARD/SCENE_092
- SHOE/SCENE_868
- KETTLE/SCENE_399
- CLOCK/SCENE_524
- HAT/SCENE_039
- BACKPACK/SCENE_264
- SCISSOR/SCENE_958
- TRUCK/SCENE_232
- HANDBAG/SCENE_575
- PINEAPPLE/SCENE_182

**Table 9a. LPIPS↓**

| Method | dataset OO3D / split S[19] / P 1 | GSO S[19] 1 | CO3D R[7] 1 | WRGBD O_h 1 | WRGBD O_h 3 | Mip360 R[7] 1 | Mip360 R[7] 3 | DL3DV O 1 | DL3DV O 3 | T&T O 1 | T&T O 3 | T&T O 6 | T&T O 9 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SV3D [19] | 0.158 | 0.140 | - | - | - | - | - | - | - | - | - | - | - |
| DepthSplat [45] | 0.610 | 0.543 | 0.756 | 0.732 | 0.588 | 0.691 | 0.491 | 0.580 | 0.405 | 0.774 | 0.706 | 0.611 | 0.487 |
| CAT3D [8] | - | - | - | - | - | - | 0.488 | - | - | - | - | - | - |
| ViewCrafter [9] | 0.634 | 0.559 | 0.789 | 0.775 | 0.603 | 0.723 | 0.540 | 0.616 | 0.576 | 0.755 | 0.671 | 0.604 | 0.546 |
| SEVA | 0.160 | 0.137 | 0.445 | 0.423 | 0.289 | 0.573 | 0.364 | 0.484 | 0.316 | 0.571 | 0.463 | 0.387 | 0.328 |

(a) **LPIPS↓**

**Table 9b. SSIM↑**

| Method | OO3D S[19] 1 | GSO S[19] 1 | CO3D R[7] 1 | WRGBD O_h 1 | WRGBD O_h 3 | Mip360 R[7] 1 | Mip360 R[7] 3 | DL3DV O 1 | DL3DV O 3 | T&T O 1 | T&T O 3 | T&T O 6 | T&T O 9 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SV3D [19] | 0.850 | 0.880 | - | - | - | - | - | - | - | - | - | - | - |
| DepthSplat [45] | 0.549 | 0.612 | 0.385 | 0.234 | 0.335 | 0.206 | 0.291 | 0.349 | 0.452 | 0.304 | 0.315 | 0.326 | 0.367 |
| CAT3D [8] | - | - | - | - | - | - | 0.294 | - | - | - | - | - | - |
| ViewCrafter [9] | 0.463 | 0.575 | 0.277 | 0.225 | 0.321 | 0.199 | 0.264 | 0.323 | 0.400 | 0.312 | 0.328 | 0.337 | 0.343 |
| SEVA | 0.857 | 0.873 | 0.536 | 0.505 | 0.603 | 0.282 | 0.377 | 0.360 | 0.480 | 0.342 | 0.385 | 0.427 | 0.452 |

(b) **SSIM↑**

Table 9. **LPIPS↓ (top) and SSIM↑ (bottom) on large-viewpoint set NVS.** For all results with $P = 1$, we sweep the unit length for camera normalization due to the model's scale ambiguity. Underlined numbers are run by us using the officially released code.

**Table 10a. LPIPS↓**

| Method | small-viewpoint RE10K | LLFF | DTU | CO3D | large-viewpoint Mip360 |
|---|---|---|---|---|---|
| ZipNeRF [47] | 0.332 | 0.373 | 0.383 | 0.652 | 0.705 |
| ZeroNVS [13] | 0.422 | 0.512 | 0.223 | 0.566 | 0.680 |
| ReconFusion [7] | 0.144 | 0.203 | 0.124 | 0.398 | 0.585 |
| CAT3D [8] | 0.132 | 0.181 | 0.121 | 0.351 | 0.515 |
| SEVA | 0.078 | 0.164 | 0.107 | 0.256 | 0.435 |

(a) **LPIPS↓**

**Table 10b. SSIM↑**

| Method | small-viewpoint RE10K | LLFF | DTU | CO3D | large-viewpoint Mip360 |
|---|---|---|---|---|---|
| ZipNeRF [47] | 0.774 | 0.574 | 0.601 | 0.496 | 0.271 |
| ZeroNVS [13] | 0.675 | 0.359 | 0.716 | 0.581 | 0.316 |
| ReconFusion [7] | 0.910 | 0.724 | 0.875 | 0.662 | 0.358 |
| CAT3D [8] | 0.917 | 0.731 | 0.844 | 0.666 | 0.377 |
| SEVA | 0.961 | 0.735 | 0.867 | 0.702 | 0.454 |

(b) **SSIM↑**

Table 10. **LPIPS↓ (top) and SSIM↑ (bottom) on 3DGS renderings for set NVS.** Results are reported on the ReconFusion [7] split with $P = 3$.

**Table 11a. LPIPS↓**

| Method | split small-viewpoint V[9] RE | CO3D | T&T | large-viewpoint O RE | DTU | WR | DL | T&T |
|---|---|---|---|---|---|---|---|---|
| MotionCtrl [46] | 0.386 | 0.443 | 0.473 | - | - | - | - | - |
| DepthSplat [45] | 0.224 | 0.532 | 0.415 | 0.134 | 0.253 | 0.452 | 0.572 | 0.685 |
| ViewCrafter [9] | 0.178 | 0.283 | 0.324 | 0.120 | 0.187 | 0.346 | 0.566 | 0.674 |
| SEVA | 0.231 | 0.318 | 0.353 | 0.079 | 0.159 | 0.284 | 0.329 | 0.514 |
| SEVA (+ temp.) | 0.228 | 0.312 | 0.356 | 0.078 | 0.156 | 0.280 | 0.328 | 0.510 |

(a) **LPIPS↓**

**Table 11b. SSIM↑**

| Method | split small-viewpoint V[9] RE | CO3D | T&T | large-viewpoint O RE | DTU | WR | DL | T&T |
|---|---|---|---|---|---|---|---|---|
| MotionCtrl [46] | 0.587 | 0.502 | 0.384 | - | - | - | - | - |
| DepthSplat [45] | 0.723 | 0.486 | 0.408 | 0.844 | 0.723 | 0.447 | 0.539 | 0.497 |
| ViewCrafter [9] | 0.798 | 0.641 | 0.563 | 0.868 | 0.739 | 0.464 | 0.523 | 0.456 |
| SEVA | 0.693 | 0.585 | 0.437 | 0.890 | 0.756 | 0.613 | 0.475 | 0.363 |
| SEVA (+ temp.) | 0.695 | 0.590 | 0.436 | 0.891 | 0.760 | 0.616 | 0.476 | 0.369 |

(b) **SSIM↑**

Table 11. **LPIPS↓ (top) and SSIM↑ (bottom) on trajectory NVS.** For the V [9] split, $P = 1$ with unit length swept; for the O split, $P = 3$. RE, WR, and DL denote RE10K, WRGBD, and DL3DV, respectively. Underlined numbers are run by us using the officially released code.

- TRAIN/SCENE_033
- REMOTE_CONTROL/SCENE_453
- BOWL/SCENE_673
- TV/SCENE_062

Full test scenes are chosen for the remaining datasets.

**Choice of input and target views.** We follow the same setup for splits adopted from previous works, by using the same set of input and target views. For split defined ourselves, we detail the choice of views as below. For the WildRGB-D [40] dataset, which consists of scenes captured while orbiting around an object, we define two splits with different difficulty levels. $O_e$ represents the easy set, where each scene is trimmed to one-third of the original sequence (*i.e.*, approximately 120 degrees of rotation). In contrast, $O_h$ corresponds to the hard set, using the full original sequence (*i.e.*, approximately 360 degrees of rotation). We first uniformly subsample 21 frames from the scene, and randomly choose $P$ frames as input views with the remaining frames as target views. For each scene from DL3DV-140 [43] and Tanks and Temples [44] datasets, we selected target frames by using every $8^{\text{th}}$ frame of the original sequence. For the remaining frames, we applied $K$-means clustering ($K = 32$) on a 6-dimensional vector formed by concatenating the camera translation and the unit vector of the camera direction.
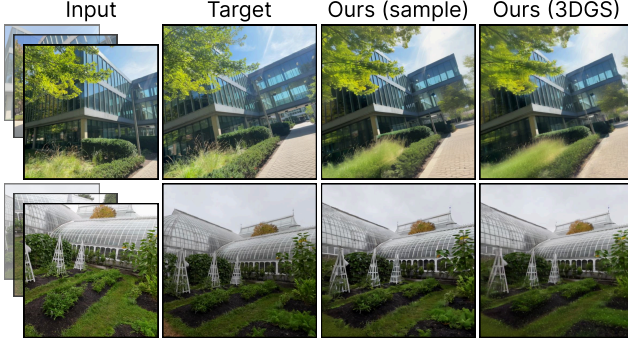
Figure 13. **3DGS versus samples.** The model generates consistent renderings that closely resemble those from 3DGS [3], with minimal perceptual differences.



Figure 14. **Padding.** Padding the last elements within one forward reduces artifacts compared to changing $T$.

## D. Additional Experiments

### D.1. Qualitative Results

We provide additional single-view conditioning sampling results with a diverse set of camera motions and effects on a variety of image prompts: a text-prompted object-centric scene (Fig. 15), a text-prompted scene (Fig. 16), a real-world object-centric scene (Fig. 17), and a real-world scene (Fig. 18). SEVA demonstrates strong generalization, adapting robustly across a wide range of scenarios.

### D.2. Quantitative Results

We provide additional quantitative evaluation results of our model against baselines on set NVS and trajectory NVS, measured using LPIPS [48] and SSIM [49], in Tabs. 8 to 10 and Appendix C.

### D.3. Discussion

**Samples *versus* 3DGS.** We compare our samples to their 3DGS distillation on the O split of DL3DV, shown in Fig. 13. First, we note that our samples contain plausible hallucinations when uncertainty is high (first row, building on the right). Second, we note that our 3DGS renderings remain sharp and are close to the samples. These results suggest that our samples are 3D consistent enough.

**Padding $T$ when $P + Q < T$.** We analyze the effect of different padding strategies when $P + Q < T$ in Fig. 14. We observe that zero-shot generalization of $T$ to $P + Q$ without padding leads to abnormal color overflows. This is in stark contrast to the excessive blurriness observed when generalizing $T$ when $P + Q >> T$ in sparse-view regime (Sec. 4.5). Hypothetically, sampling with a $T$ unseen during training induces a distribution shift in the attention scores [70]. Specifically, a smaller $T$ sharpens the attention distribution, whereas a larger $T$ disperses it. This shift may explain the contrasting behavior observed when
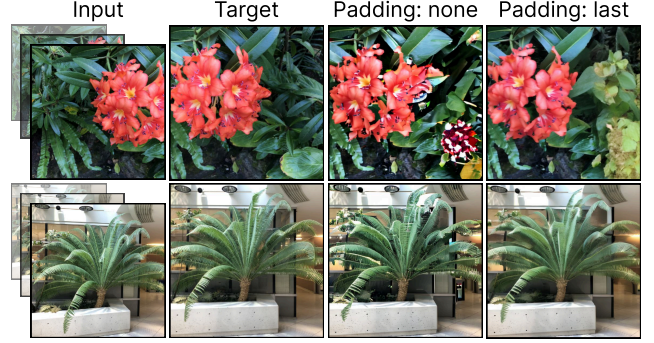
using the model for sampling. Training the model with a dynamically varying $T$ during training could mitigate this issue by exposing the model to a broader range of attention score distributions, improving generalization across different $T$.

**Artifacts on long-trajectory NVS.** We observe that the results tend to become increasingly saturated, particularly when the target views are far from the input views and share no content overlap, such as in open-ended exploration and navigation. The concurrent work [71] explores the concept of Diffusion Forcing [72] for long video rollouts, achieving high-generation quality. Applying diverse noise to the input views during training can be beneficial, as it enables the refinement of high-level details in all anchor views within the memory bank during sampling, thereby mitigating the accumulation of saturation. We leave this for future work.

Figure 15. **Diverse camera motions and effects.** Single-view conditioning with a text-prompted object-centric scene. The image is generated using SD 3.5 [73] with the text prompt, "*A cute firefly dragon in its natural habitat.*"
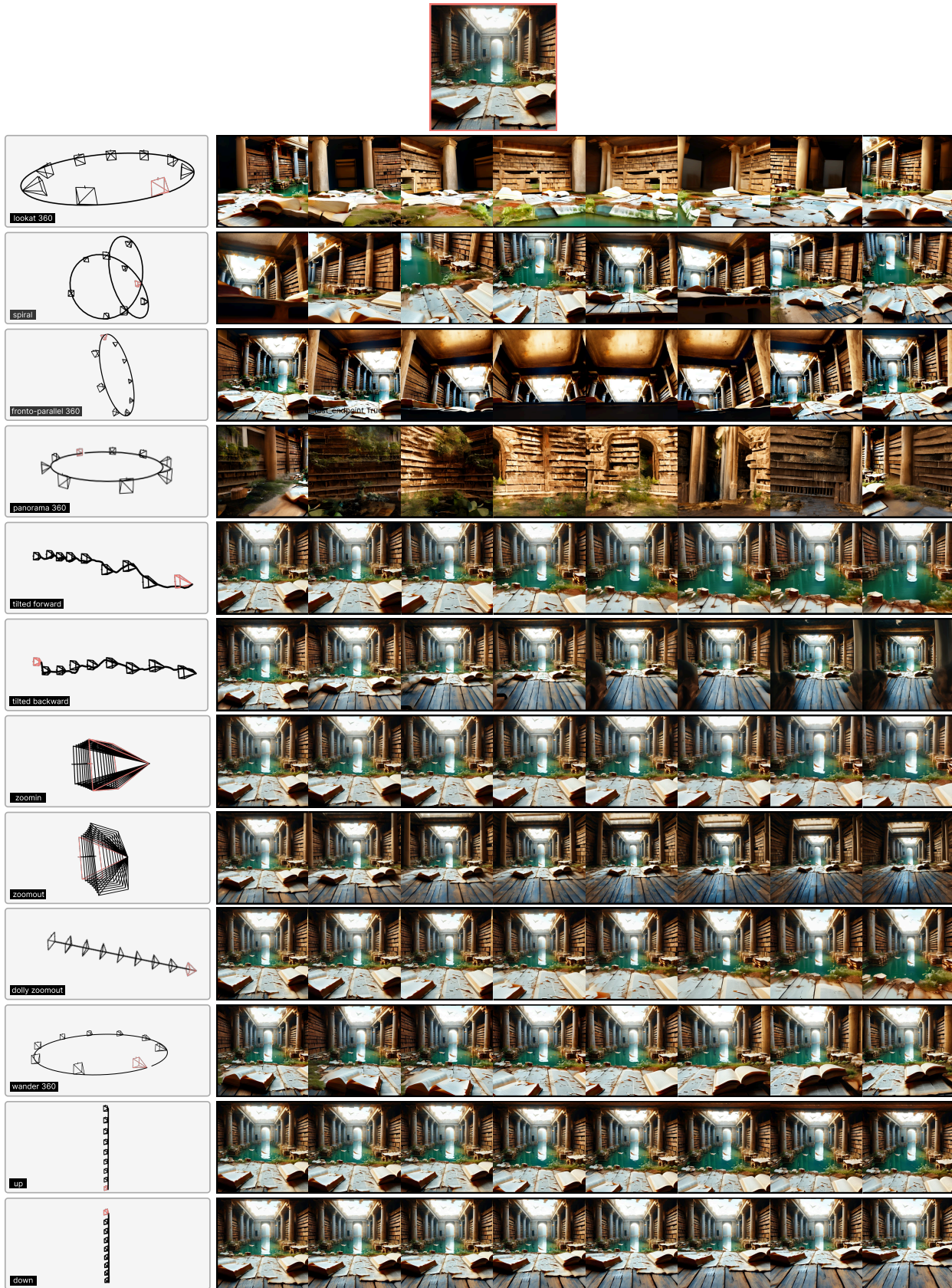
Figure 16. **Diverse camera motions and effects.** Single-view conditioning with a text-prompted scene. The image is generated using SD 3.5 [73] with the text prompt, "*Wide view of the interior of the famed Library of Alexandria, elegantly set behind a time-worn wreckage by a lake, hinting at the relentless passage of time. The surroundings are lit by the light of a late afternoon sun, gently cast, immersing the area in a sentimental luminescence.*"

Figure 17. **Diverse camera motions and effects.** Single-view conditioning with a real-life object-centric scene.
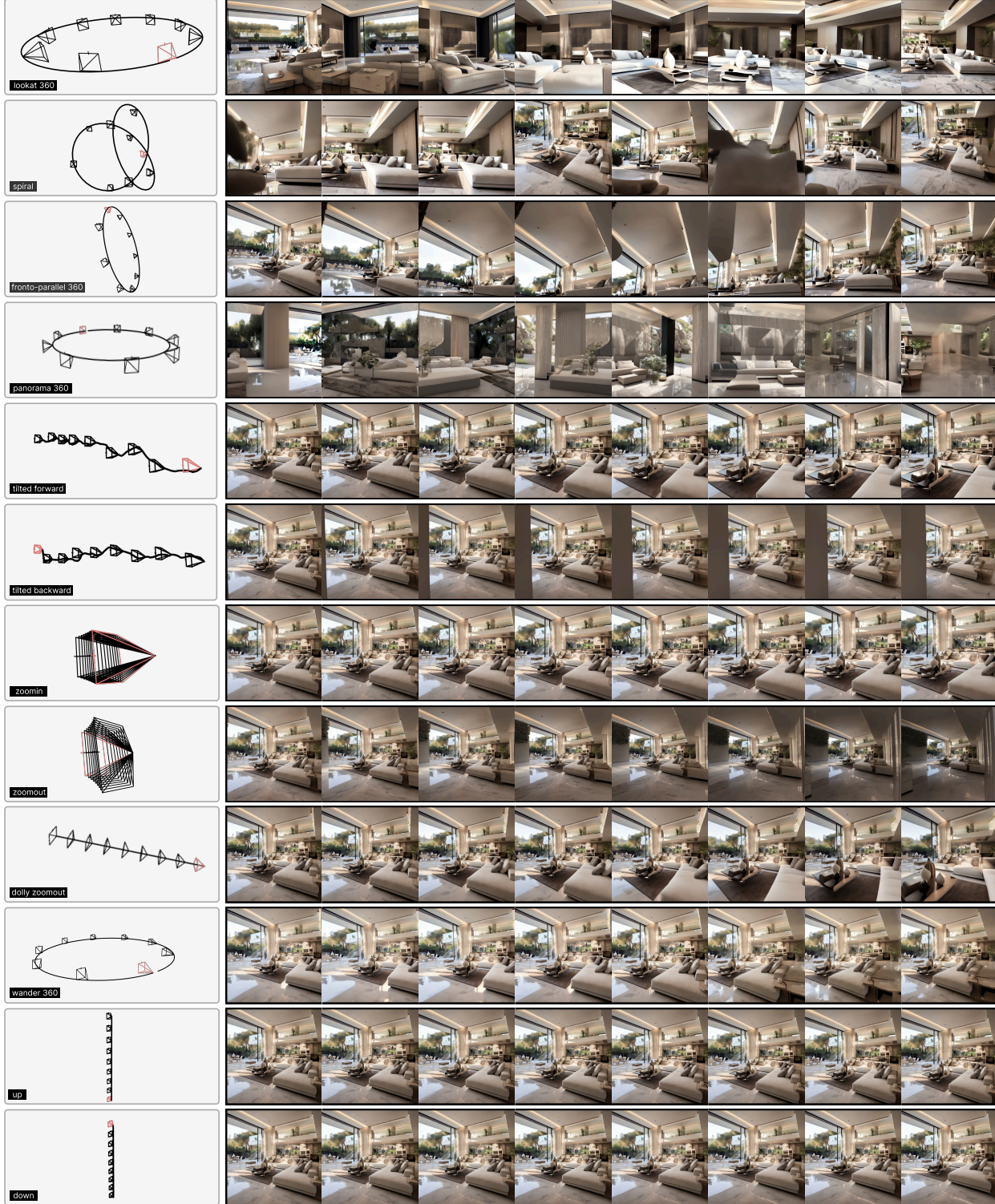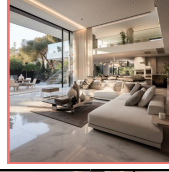
Figure 18. **Diverse camera motions and effects.** Single-view conditioning with a real-life scene.