# SpaceTimePilot: Generative Rendering of Dynamic Scenes Across Space and Time

Zhening Huang[1,2]   Hyeonho Jeong[2]   Xuelin Chen[2]   Yulia Gryaditskaya[2]

Tuanfeng Y. Wang[2]   Joan Lasenby[1]   Chun-Hao Huang[2]

[1]University of Cambridge   [2]Adobe Research

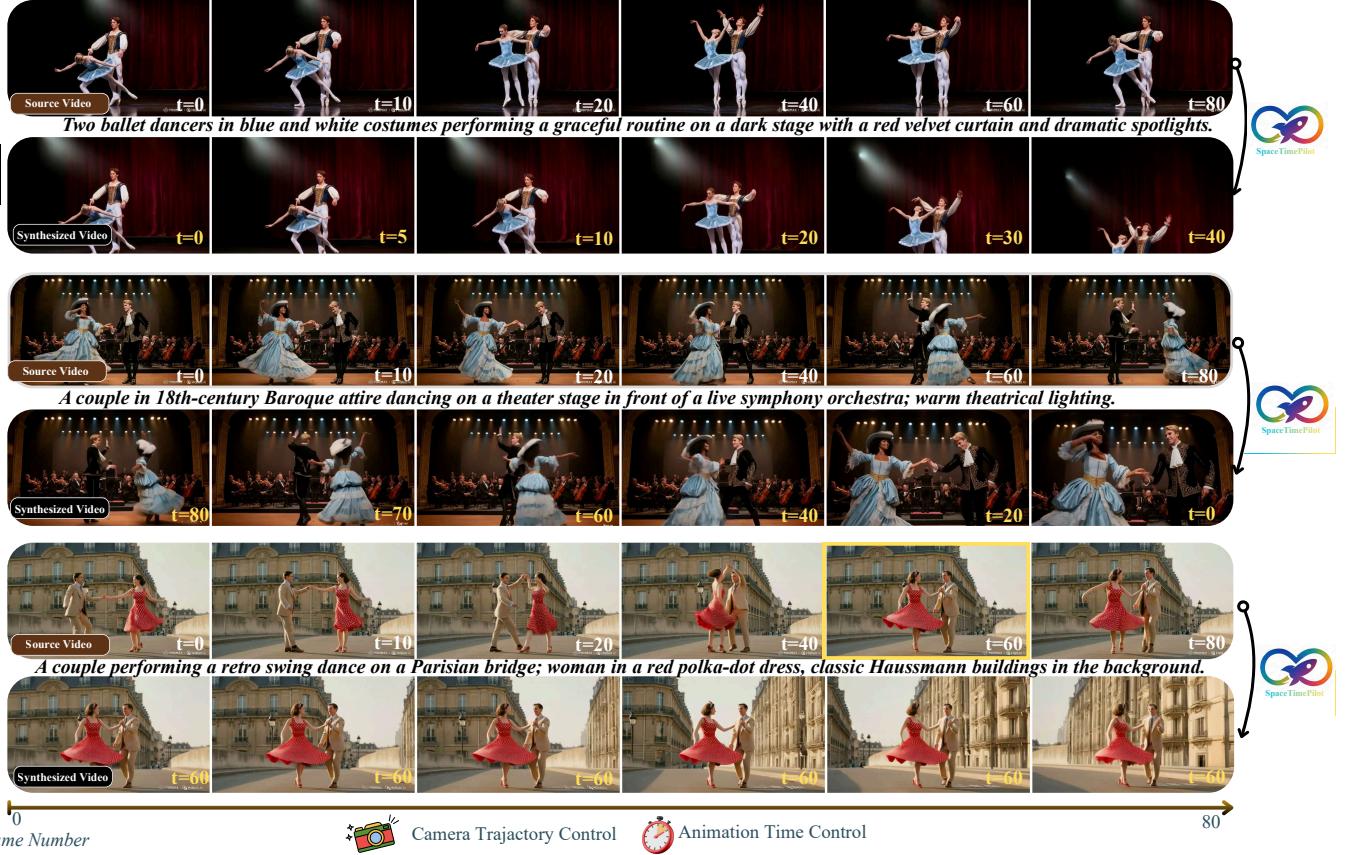https://zheninghuang.github.io/Space-Time-Pilot/



Figure 1. **SpaceTimePilot** enables unified control over both camera and time within a single diffusion model, producing continuous and coherent videos along arbitrary space–time trajectories. Given a source video (odd rows), our model synthesizes new videos (even rows) with retimed motion sequences, including slow motion, reverse motion, and bullet time, while precisely controlling camera movement according to a given camera trajectory.

## Abstract

We present **SpaceTimePilot**, a video diffusion model that disentangles space and time for controllable generative rendering. Given a monocular video, **SpaceTimePilot** can independently alter the camera viewpoint and the motion sequence within the generative process, re-rendering the scene for continuous and arbitrary exploration across space and time. To achieve this, we introduce an effective animation time-embedding mechanism in the diffusion process, allowing explicit control of the output video's motion sequence with respect to that of the source video. As no datasets provide paired videos of the same dynamic scene with continuous temporal variations, we propose a simple

1

*yet effective **temporal-warping training scheme** that repurposes existing multi-view datasets to mimic temporal differences. This strategy effectively supervises the model to learn temporal control and achieve robust space–time disentanglement. To further enhance the precision of dual control, we introduce two additional components: an improved camera-conditioning mechanism that allows altering the camera from the first frame, and Cam×Time, the first synthetic Space and Time full-coverage rendering dataset that provides fully free space–time video trajectories within a scene. Joint training on the temporal-warping scheme and the Cam×Time dataset yields more precise temporal control. We evaluate **SpaceTimePilot** on both real-world and synthetic data, demonstrating clear space–time disentanglement and strong results compared to prior work.*

## 1. Introduction

Videos are 2D projections of an evolving 3D world, where the underlying generative factors consist of spatial variation (camera viewpoint) and temporal evolution (dynamic scene motion). Learning to understand and disentangle these factors from observed videos is fundamental for tasks such as scene understanding, 4D reconstruction, video editing, and generative rendering, to name a few. In this work, we approach this challenge from the perspective of generative rendering. Given a single observed video of a dynamic scene, our goal is to synthesize novel views (reframe/reangle) and/or at different moments in time (retime), while remaining faithful to the underlying scene dynamic.

A common strategy is to first reconstruct dynamic 3D content from 2D observations, *i.e.*, perform 4D reconstruction, and then re-render the scene. These methods model both spatial and temporal variations using representations such as NeRFs [22, 25] or Dynamic Gaussian Splatting [15, 42], often aided by cues like geometry [27, 28], optical flow [19, 20], depth [6, 47], or long-term 2D tracks [17, 37]. However, even full 4D reconstructions typically show artifacts under novel viewpoints. More recent work [21, 43] uses multi-view video diffusion to generate sparse, time-conditioned views and refines them via Gaussian-splatting optimization, but rendering quality remains limited. Advances in video diffusion models [3–5, 12, 16, 30, 39, 46, 51] further enable camera re-posing with more lightweight point cloud representations, reducing the need for heavy 4D reconstruction. While effective in preserving identity, their reliance on per-frame depth and reprojection limits robustness under large viewpoint changes. To mitigate this, newer approaches condition generation solely on camera parameters, achieving strong novel-view synthesis on both static [14] and dynamic scenes [2, 9, 33]. Autoregressive models like Genie-3 [29] even enable interactive scene exploration from a single image, showing that
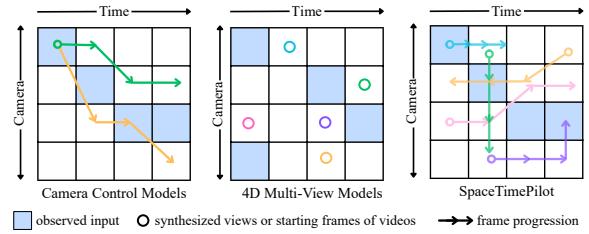


Figure 2. **Space–time controllability across methods.** Blue cells denote the input video/views, while arrows and dots indicate generated continuous videos or sparse frames. Camera-control V2V models [2, 33] modify only the camera trajectory while keeping time strictly monotonic. 4D multi-view models [21, 43] synthesize discrete sparse views conditioned on space and time, but do not generate continuous video sequences. *SpaceTimePilot* enables free movement along both the camera and time axes with full control over direction and speed, supporting bullet-time, slow-motion, reverse playback, and mixed space–time trajectories.

diffusion models can encode implicit 4D priors. Nonetheless, despite progress in spatial viewpoint control, current methods still lack full 4D exploration, *i.e.*, the ability to navigate scenes freely across both space and time.

In this work, we introduce *SpaceTimePilot*, the first video diffusion model that enables joint spatial and temporal control. SpaceTimePilot introduces a new notion of "animation time" to capture the *temporal status* of scene dynamics in the source video. As such, it naturally disentangles temporal control and camera control by expressing them as two independent signals. A high-level comparison between our approach and prior methods is illustrated in Fig. 2. Unlike previous methods, SpaceTimePilot enables free navigation along both the camera and time axes. Training such a model requires dynamic videos that exhibit multiple forms of temporal playback while simultaneously being captured under multiple camera motions, which is only feasible in a controlled studio setups. Although temporal diversity can be increased by combining multiple real datasets *e.g.* [23, 53], as done in [41, 43], this approach remains suboptimal, as the coverage of temporal variation is still insufficient to learn the underlying meaning of temporal control. Existing synthetic datasets [1, 2] also do not exhibit such properties.

To address this limitation, we introduce a simple yet effective *temporal-warping training scheme* that augments existing multi-view video datasets [1, 2] to simulate diverse conditioning types while preserving continuous video structure. By warping input sequences in time, the model is exposed to varied temporal behaviors without requiring additional data collection. This simple yet crucial strategy allows the model to learn temporal control signals, enabling it to directly exhibit space–time disentanglement effects during generation. We further ablate various temporal-conditioning schemes and introduce a convolution-based

temporal-control mechanism that enables finer-grained manipulation of temporal behavior and supports effects such as bullet-time at any timestep within the video. While temporal warping increases temporal diversity, it can still entangle camera and scene dynamics – for example, temporal manipulation may inadvertently affect camera behavior. To further strengthen disentanglement, we introduce a new dataset that spans the full grid of camera–time combinations along a trajectory. Our synthetic *Cam×Time* dataset contains 180k videos rendered from 500 animations across 100 scenes and three camera paths. Each path provides full-motion sequences for every camera pose, yielding dense multi-view and full-temporal coverage. This rich supervision enables effective disentanglement of spatial and temporal control.

Experimental results show that SpaceTimePilot successfully disentangles space and time in generative rendering from single videos, outperforming adapted state-of-the-art baselines by a significant margin. Our main contributions are summarized as follows:

- We introduce SpaceTimePilot, the first video diffusion model that disentangles spatial and temporal factors to enable continuous and controllable novel view synthesis as well as temporal control from a single video.
- We propose the *temporal-warping strategy* that repurposes multi-view datasets to simulate diverse temporal variations. By training on these warped sequences, the model effectively learns temporal control without the need for explicitly constructed video pairs captured under different temporal settings.
- We propose a more precise camera–time conditioning mechanism, illustrating how viewpoint and temporal embeddings can be jointly integrated into diffusion models to achieve fine-grained spatiotemporal control.
- We construct the Cam×Time Dataset, providing dense spatiotemporal sampling of dynamic scenes across camera trajectories and motion sequences. This dataset supplies the necessary supervision for learning disentangled 4D representations and supports precise camera–time control in generative rendering.

## 2. Related work

We aim to re-render a video from new viewpoints with temporal control, a task closely related to Novel View Synthesis (NVS) from monocular video inputs.

**Video-based NVS.** Prior video-based NVS methods can be broadly characterized along two axes: (i) whether they target static or dynamic scenes, and (ii) whether they incorporate explicit 3D geometry in the generation pipeline.

For static scenes, geometry-based methods reconstruct scene geometry from the input frames and use diffusion models to complete or hallucinate regions that are unseen under new viewpoints [13, 31, 44, 49]. Although these approaches achieve high rendering quality, they rely on heavy 3D preprocessing. Geometry-free approaches [2, 33, 52] bypass explicit geometry and directly condition the diffusion process on observed views and camera poses to synthesize new viewpoints.

For dynamic scenes, inpainting-based methods such as TrajectoryCrafter [48], ReCapture [50], and Reangle [13] also adopt warp-and-inpaint pipelines, while GEN3C [31] extends this with an evolving 3D cache and EPiC [40] improves efficiency via a lightweight ControlNet framework. Geometry-free dynamic models [1, 2, 33, 35, 36] instead learn camera-conditioned generation from multi-view or 4D datasets (*e.g.*, Kubric-4D [7]), enabling smoother and more stable NVS with minimal 3D inductive bias. Proprietary systems like Genie 3 [29] further demonstrate real-time, continuous camera control in dynamic scenes, underscoring the potential of video diffusion models for interactive viewpoint manipulation.

**Disentangling Space and Time.** Despite great progress in camera controllability (space), the methods discussed above do not address temporal control (time). Meanwhile, disentangling spatial and temporal factors has become a central focus in 4D scene generation, recently advanced through diffusion-based models. 4DiM [41] introduces a Masked FiLM mechanism that defaults to identity transformations when conditioning signals (e.g., camera pose or time) are absent, enabling unified representations across both static and dynamic data through multi-modal supervision. Similarly, CAT4D [43] leverages multi-view images to conduct 4D dynamic reconstruction to achieve space–time disentanglement but remains constrained by its reliance on explicit 4D reconstruction pipelines, which limits scalability and controllability. In contrast, our approach builds upon text-to-video diffusion models and introduces a new temporal embeddings module and refined camera conditioning to achieve fully controllable 4D generative reconstruction.

## 3. Method

We introduce SpaceTimePilot, a method that takes a source video $V_{src} \in \mathbb{R}^{F \times C \times H \times W}$ as input and synthesizes a target video $V_{trg} \in \mathbb{R}^{F \times C \times H \times W}$, following an input camera trajectory $\mathbf{c}_{trg} \in \mathbb{R}^{F \times 3 \times 4}$ and temporal control signal $\mathbf{t}_{trg} \in \mathbb{R}^F$. Here, $F$ denotes the number of frames, $C$ the number of color channels, and $H$ and $W$ are the frame height and width, respectively. Each $\mathbf{c}_{trg}^f \in \mathbb{R}^{3 \times 4}$ represents the camera extrinsic parameters (rotation and translation) at frame $f$, with respect to the 1st frame of $V_{src}$. The target video $V_{trg}$ preserves the scene's underlying dynamics, geometry, and appearance in $V_{src}$, while adhering to the camera motion and temporal progression specified by $\mathbf{c}_{trg}$ and $\mathbf{t}_{trg}$. A key feature of our method is the disentanglement

of spatial and temporal factors in the generative process, enabling effects such as bullet-time and retimed playback from novel viewpoints (see Fig. 1).

## 3.1. Preliminaries

Our framework builds upon recent advances in large-scale text-to-video diffusion models and camera-conditioned video generation. We adopt a latent video diffusion backbone similar to modern text-to-video foundation models [34], consisting of a 3D Variational Auto-Encoder (VAE) for latent compression and a Transformer-based denoising model (DiT) operating over multi-modal tokens.

Additionally, our design draws inspiration from ReCamMaster [2], which introduces explicit camera conditioning for video synthesis. Given an input camera trajectory $\mathbf{c} \in \mathbb{R}^{F \times 3 \times 4}$, spatial conditioning is achieved by first projecting the camera sequence to the space of video tokens and adding it to the features:

$$x' = x + \mathcal{E}_{\text{cam}}\left(\mathbf{c}\right), \qquad (1)$$

where $x$ is the output of the patchifying module and $x'$ is the input to self-attention layers. The camera encoder $\mathcal{E}_{\text{cam}}$ maps each flattened $3 \times 4$ camera matrix (12-dimensional) into the target feature space, while also transforming the temporal dimension from $F$ to $F'$.

## 3.2. Disentangling Space and Time

We achieve spatial and temporal disentanglement through a two-fold approach: a dedicated time representation and specialized datasets.

### 3.2.1. Time representation

Recent video diffusion models include position embeddings for latent frame index $f'$, such as RoPE($f'$). However, we found using RoPE($f'$) for temporal control to be ineffective, as it interferes with camera signals: RoPE($f'$) often constrains both temporal and camera motion simultaneously. To address space and time disentanglement, we introduce a dedicated time control parameter $\mathbf{t} \in \mathbb{R}^F$. By manipulating $\mathbf{t}_{\text{trg}}$, we can control the temporal progression of the synthesized video $V_{\text{trg}}$. For example, setting $\mathbf{t}_{\text{trg}}$ to a constant locks $V_{\text{trg}}$ to a specific timestamp in $V_{\text{src}}$, while reversing the frame indices produces a playback of $V_{\text{src}}$ in reverse.

(Top) For multi-view dynamic scene datasets, a set of temporal warping operations, including reverse, playback, zigzag motion, slow motion, and freeze are apppplied with teh source video as standford. This gives explicit supervision for temporal control, without constructing additional temporally varied training data.

(Bottom) Existing camera-control and joint dataset training rely on monotonic time progression and static scene videos, making it difficult for models to understand temporal variation. The introduced temporal mappings from
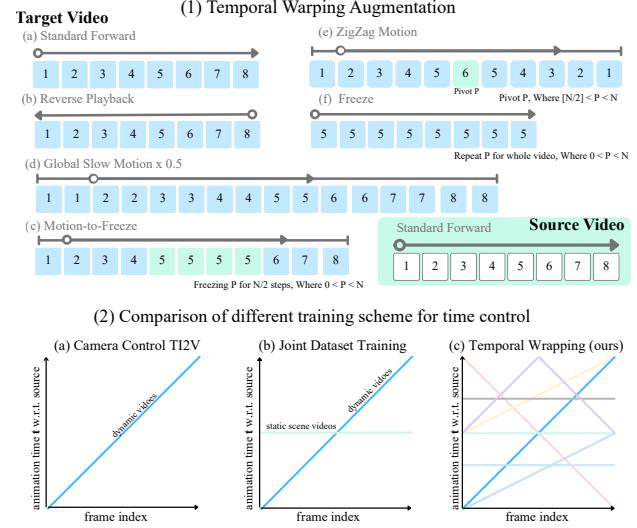


Figure 3. **Temporal Wrapping for Spatiotemporal Disentanglement.** (Top) For multi-view dynamic scene datasets [2], a set of temporal warping operations (e.g. reverse playback, zigzag motion, slow motion, and freeze) are applied to the target video, with the source video kept as the standard forward reference, providing explicit supervision for temporal control . (Bottom) Compared with existing camera-control [2, 33] and joint-dataset training strategies [41, 43], which rely on monotonic time progression and static-scene videos to demonstrate temporal differences, Temporal Wrapping provide much more diverse and explicit signals of temporal variation, leading to disentanglement of space and time.

multi-view video data, which provide diverse and clear signal on tempral variation, and directly lead to disentanglement of space and time.

**Time Embedding.** To inject temporal control into the diffusion model, we analyze several approaches. First, we can encode time similar to a frame index using RoPE embedding. However, we find it less suitable for time control (visual evaluations are provided in Supp. Mat.). Instead, we adopt sinusoidal time embeddings applied at the latent frame $f'$ level, which provide a stable and continuous representation of each frame's temporal position and offer a favorable trade-off between precision and stability. We further observe that each latent frame corresponds to a continuous temporal chunk, and propose using embeddings of original frame indices $f$ to support finer granularity of time control. To accomplish this, we introduce a time encoding approach $\mathcal{E}_{\text{ani}}(\mathbf{t})$, where $\mathbf{t} \in \mathbb{R}^F$. We first compute the sinusoidal time embeddings to represent the temporal sequence, $\mathbf{e}_{\text{src}} = \text{SinPE}(\mathbf{t}_{\text{src}})$, $\mathbf{e}_{\text{trg}} = \text{SinPE}(\mathbf{t}_{\text{trg}})$, where $\mathbf{t}_{\text{src}}, \mathbf{t}_{\text{trg}} \in \mathbb{R}^F$. Next, we apply two 1D convolution layers to progressively project these embeddings into the latent frame space, $\widetilde{\mathbf{e}} = \text{Conv1D}_2(\text{Conv1D}_1(\mathbf{e}))$. Finally, we add these time features to the camera features and video tokens

embeddings, updating Eq. (1) as follows:

$$x' = x + \mathcal{E}_{\text{cam}}(\mathbf{c}) + \mathcal{E}_{\text{ani}}(\mathbf{t}). \quad (2)$$

In Sec. 4.2, we compare our approach with alternative conditioning strategies, such as using sinusoidal embeddings where $\mathbf{t}_{\text{src}}, \mathbf{t}_{\text{trg}}$ are directly defined in $\mathbb{R}^{F'}$, and employing an MLP instead of a 1D convolution for compression. We demonstrate both qualitatively and quantitatively the advantages of our proposed method.

### 3.2.2. Datasets

To enable temporal manipulation in our approach, we require paired training data that includes examples of time remapping. Achieving spatial-temporal disentanglement further requires data containing examples of both camera and temporal controls. To the best of our knowledge, no publicly available datasets satisfy these requirements. Only a few prior works, such as 4DiM [41] and CAT4D [43], have attempted to address spatial-temporal disentanglement. A common strategy is to jointly train on static-scene datasets and multi-view video datasets [23, 53]. The limited control variability in these datasets leads to confusion between temporal evolution and spatial movement, resulting in entangled or unstable behaviors [41, 43]. We address this limitation by augmenting existing multi-view video data with temporal warping and by proposing a new synthetic dataset.

**Temporal Warping Augmentation.** We introduce simple augmentations that add controllable temporal variations to multi-view video datasets. During training, given a source video $V_{\text{src}} = \{I_{\text{src}}^f\}_{f=1}^F$ and a target video $V_{\text{trg}} = \{I_{\text{trg}}^f\}_{f=1}^F$, we apply a temporal warping function $\tau : [1, F] \to [1, F]$ to the target sequence, producing a warped video $V'_{\text{trg}} = \{I_{\text{trg}}^{\tau(f)}\}_{f=1}^F$. The source animation timestamps are uniformly sampled, $\mathbf{t}_{\text{src}} = 1 : F$. Warped timestamps, $\mathbf{t}_{\text{trg}} = \tau(\mathbf{t}_{\text{src}})$, introduce non-linear temporal effects (see Fig. 3 top b–e): (i) reversal, (ii) acceleration, (iii) freezing, (iv) segmental slow motion, and (v) zigzag motion, in which the animation repeatedly reverses direction. After these augmentations, the paired video sequences $(V_{\text{src}}, V'_{\text{trg}})$ differ in both camera trajectories and temporal dynamics, providing the model with a clear signal for learning disentangled spatiotemporal representations.

**Synthetic Cam×Time Dataset for Precise Spatiotemporal Control.** While our temporal warping augmentations encourage strong disentanglement between spatial and temporal factors, achieving fine-grained and continuous control — that is, smooth and precise adjustment of temporal dynamics — benefits from a dataset that systematically covers both dimensions. To this end, we construct *Cam×Time*, a new synthetic spatiotemporal dataset rendered in Blender. Given a camera trajectory and an animated subject, Cam×Time exhaustively samples the camera–time grid, capturing each dynamic scene across diverse

Table 1. **Comparison of existing multi-view datasets for camera and temporal control against Cam×Time.** Cam×Time provides full-grid rendering (Figure 4), enabling target videos to sample arbitrary temporal variations over the full range from 0 to 120.

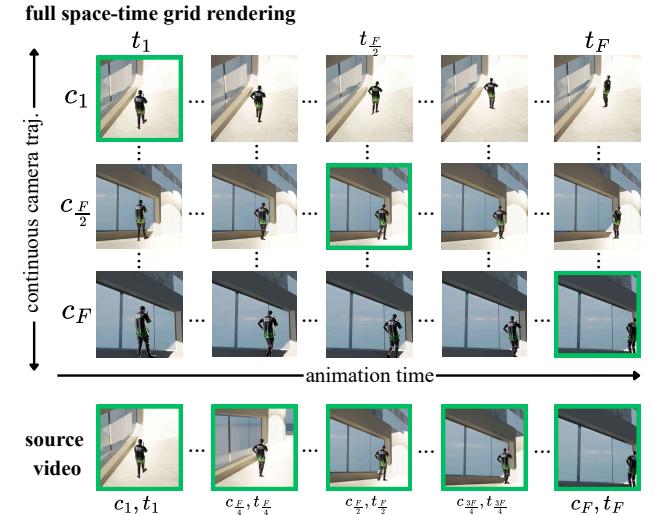| Dataset | Dynamic scenes | Src. Time: $t_{\text{src}}$ | Tgt. Time: $t_{\text{trg}}$ | Camera |
|---|---|---|---|---|
| RE10k [53] | ✗ | 1 | 1 | Moving |
| DL3DV10k [23] | ✗ | 1 | 1 | Moving |
| MannequinChallenge [32] | ✗ | 1 | 1 | Moving |
| Kubric-4D [33] | ✔ | 1:60 | 1:60 | Moving |
| ReCamMaster [2] | ✔ | 1:80 | 1:80 | Moving |
| SynCamMaster [1] | ✔ | 1:80 | 1:80 | Fixed |
| **Cam×Time (ours)** | ✔ | 1:120 | $\{1, 2, \ldots, 120\}^{120}$ | Moving |

**full space-time grid rendering**



Figure 4. *Cam×Time* **dataset visualization**. (Top) A space-time grid defined by a camera trajectory $\mathbf{c} = [c_1, ..., c_F]$ and animation status $\mathbf{t} = [t_1, ..., t_F]$. Cam×Time renders images for all $(c, t)$ pairs, covering the full grid for learning disentangled spatial and temporal control. Any two sampled sequences of $F$ frames from the grid can form a source-target pair. (Bottom) One typical choice of source videos is taking the diagonal cells in green.

combinations of camera viewpoints and temporal states $(\mathbf{c}, \mathbf{t})$, as illustrated in Fig. 4. The source video is obtained by sampling the diagonal frames of the dense grid (Fig. 4 (bottom)), while the target videos are obtained by more free-form sampling of continuous sequences. We compare Cam×Time against existing datasets in Tab. 1. While [23, 32, 53] are real videos with complex camera path annotations, they either do not provide time-synchronized video pairs [32] or only provide pairs of static scenes [23, 53]. Synthetic multi-view video datasets [1, 2, 33] provide pairs of dynamic videos but do not allow training for time control. In contrast, Cam×Time enables fine-grained manipulation of both camera motion and temporal dynamics, enabling bullet-time effects, motion stabilization, and flexible combinations of the controls. We designate part of Cam×Time as a test set, aiming for it to serve as a benchmark for controllable video generation. We will release it to support future research on fine-grained spatiotemporal modeling.

5

Table 2. Quantitative comparison across temporal controls (*Direction (forward, backward motion)*, *Speed (slow modes)*, *Bullet Time*). We report PSNR↑, SSIM↑, and LPIPS↓. Best results are in **bold**. SpaceTimeMethod showcase best performance for temporal control overall.

| Method | PSNR↑ | | | | SSIM↑ | | | | LPIPS↓ | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Dir. | Speed | Bullet | Avg | Dir. | Speed | Bullet | Avg | Dir. | Speed | Bullet | Avg |
| ReCamM+preshuffled[†] | 17.13 | 14.84 | 14.61 | 15.52 | 0.6623 | 0.6050 | 0.5965 | 0.6213 | 0.3930 | 0.4793 | 0.4863 | 0.4529 |
| ReCamM+jointdata | 18.32 | 17.57 | 17.69 | 17.86 | 0.7322 | 0.7220 | 0.7209 | 0.7250 | 0.2972 | 0.3158 | 0.3089 | 0.3073 |
| **SpaceTimePilot (Ours)** | **21.75** | **20.87** | **20.85** | **21.16** | **0.7725** | **0.7645** | **0.7653** | **0.7674** | **0.1697** | **0.1917** | **0.1677** | **0.1764** |

[†] Uses simple frame-rearrangement operators (reversal, repetition, freezing) applied prior to inference to emulate temporal manipulation.

### 3.3. Precise Camera Conditioning

We aim for full camera trajectory control in the target video. In contrast, the previous novel-view synthesis approach [2] assumes that the first frame is identical in source and target videos and that the target camera trajectory is defined relative to it. This stems from the two limitations. First, the existing approach ignores the source video trajectory, yielding suboptimal source features computed using the target trajectory for consistency:

$$x'_{\text{src}} = x_{\text{src}} + \mathcal{E}_{\text{cam}}\left(\mathbf{c}_{\text{trg}}\right), \quad x'_{\text{trg}} = x_{\text{trg}} + \mathcal{E}_{\text{cam}}\left(\mathbf{c}_{\text{trg}}\right).$$

Second, it is trained on datasets where the first frame is always identical across the source and target videos. This latter limitation is addressed in our training datasets design.

To overcome the former, we devise a *source-aware camera conditioning*. We estimate camera poses for both the source and target videos using a pretrained pose estimator, and inject them jointly into the diffusion model to provide explicit geometric context. Eq. 2 is therefore extended into:

$$x'_{\text{src}} = x_{\text{src}} + \mathcal{E}_{\text{cam}}\left(\mathbf{c}_{\text{src}}\right) + \mathcal{E}_{\text{ani}}\left(\mathbf{t}_{\text{src}}\right), \qquad (3)$$
$$x'_{\text{trg}} = x_{\text{trg}} + \mathcal{E}_{\text{cam}}\left(\mathbf{c}_{\text{trg}}\right) + \mathcal{E}_{\text{ani}}\left(\mathbf{t}_{\text{trg}}\right),$$
$$x' = [x'_{\text{trg}}, x'_{\text{src}}]_{\text{frame-dim}},$$

where $x'$ denotes the input of the DiT model, which is the concatenation of target and source tokens along the frame dimension. This formulation provides the model with both source and target camera context, enabling spatially consistent generation and precise control over camera trajectories.

### 3.4. Support for Longer Video Segments

Finally, to showcase the full potential of our camera and temporal control, we adopt a simple autoregressive video generation strategy, generating each new segment $V_{\text{trg}}$ conditioned on the previously generated segment $V_{\text{prv}}$ and a source video $V_{\text{src}}$ to produce longer videos.

To enable this capability during inference, we need to extend our training scenario to support conditioning on two videos, where one serves as $V_{\text{src}}$ and the other as $V_{\text{prv}}$. The source video $V_{\text{src}}$ is taken directly from the multi-view datasets or from our synthetic dataset, as was described previously. $V_{\text{prv}}$ is constructed in a similar way to $V_{\text{trg}}$ — either

using temporal warping augmentations or by sampling from the dense space-time grid of our synthetic dataset. When temporal warping is applied, $V_{\text{prv}}$ and $V_{\text{trg}}$ may originate from the same or different multi-view sequences representing the same time interval. To maintain full flexibility of control, we do not enforce any other explicit correlations between $V_{\text{prv}}$ and $V_{\text{trg}}$, apart from specifying camera parameters relative to the selected source video frame.

Note that not constraining the source and target videos to share the same first frame (as discussed in Sec. 3.3) is crucial for achieving flexible camera control in longer sequences. For instance, this design enables extended bullet-time effects: we can first generate a rotation around a selected point up to $45°$ ($V_{\text{trg,1}}$), and then continue from $45°$ to $90°$ ($V_{\text{trg,2}}$). Conditioning on two consecutive source segments allows the model to leverage information from newly generated viewpoints. In the bullet-time example, conditioning on the previously generated video enables the model to incorporate information from all newly synthesized viewpoints, rather than relying solely on the viewpoint of the corresponding moment in the source video.

## 4. Experiments

**Implementation details.** We adopt the Wan-2.1 T2V-1.3B model [34], which produces $F'=21$ latent frames and decodes them into $F=81$ RGB frames using a 3D-VAE. The network is conditioned on camera and animation-time controls as defined in Eq. 3. Unless otherwise specified, SpaceTimePilot is trained with ReCamMaster and SynCamMaster datasets with the temporal warping augmentation described in Sec. 3.2.2, along with Cam×Time. Please refer to Supp. Mat. for complete network architecture and additional training details.

### 4.1. Comparison with State-of-the-Art Baselines

#### 4.1.1. Time-Control Evaluation.

We first evaluate the retiming capability of our model. To factor out the error induced by camera control, we condition SpaceTimePilot on a fixed camera pose while varying only the temporal control signal. Experiments are performed on the withheld Cam×Time test split, which contains 50 scenes rendered with dense full-grid trajectories that can

Figure 5. **Qualitative results of SpaceTimePilot.** Our model enables fully disentangled control over camera motion and temporal dynamics. Each row shows a different combination of camera trajectory (left icons) and temporal warping (right icons). SpaceTimePilot produces coherent videos under diverse controls, including normal playback, reverse playback, bullet-time, slow-motion, replay motion, and complex camera paths (pan, tilt, zoom, and vertical motion).

be retimed into arbitrary temporal sequences. For each test case, we take a moving-camera source video but set the target camera trajectory to the first-frame pose. We then apply a range of temporal control signals, including reverse, bullet-time, zigzag, slow motion, and normal playback, to synthesize the corresponding retimed outputs. Since we have ground-truth frames for all temporal configurations, we report perceptual losses: PSNR, SSIM, and LPIPS.

We consider two baselines: (1) *ReCamM+preshuffled*: original ReCamMaster combined with input re-shuffling; and (2) *ReCamM+jointdata*: following [41, 43], we train ReCamMaster with additional static-scene datasets [18, 53] which provide only one single temporal pattern.

While frame shuffling may succeed in simple scenarios, it fails to disentangle camera and temporal control. As shown in Table 2, this approach exhibits the weakest temporal controllability. Although incorporating static-scene datasets improves performance, particularly in the bullet-time category, relying on a single temporal control pattern remains insufficient for achieving robust temporal consistency. In contrast, SpaceTimePilot consistently outperforms all baselines across all temporal configurations.

Table 3. VBench visual-quality evaluation across six dimensions. Higher is better for all metrics.

| Method | ImgQ↑ | BGCons↑ | Motion↑ | SubjCons↑ | Flicker↑ | Aesthetic↑ |
|---|---|---|---|---|---|---|
| Traj-Crafter [48] | 0.6389 | **0.9376** | 0.9888 | **0.9463** | 0.9816 | 0.5172 |
| ReCamM [2] | 0.6302 | 0.9114 | 0.9945 | 0.9181 | **0.9825** | 0.5332 |
| ReCamM+Aug | 0.6315 | 0.9165 | 0.9946 | 0.9313 | 0.9788 | 0.5385 |
| **STPilot (Ours)** | **0.6486** | **0.9199** | **0.9947** | **0.9325** | 0.9781 | **0.5315** |

### 4.1.2. Visual Quality Evaluation.

Next, we evaluate the perceptual realism of our 1800 generated videos using VBench [10]. We report all standard visual quality metrics to provide a comprehensive assessment of generative fidelity. Table 3 shows that our model achieves visual quality comparable to the baselines.

### 4.1.3. Camera-Control Evaluation.

Finlay, we evaluate the effectiveness of our camera control mechanism detailed in Sec. 3.3. Unlike the retiming evaluation above, which relies on synthetic ground-truth videos, here we construct a real-world 90-video evaluation set from OpenVideoHD [26], encompassing diverse dynamic human and object motions. Each method is evaluated across 20 camera trajectories: 10 starting from the same initial pose as the source video and 10 from different initial poses,

Table 4. Camera accuracy and first-frame estimation. For camera control, the enhanced camera control mechanism enables the generated video to start from an arbitrary camera angle while maintaining good camera accuracy.

| Method | RelRot↓ | RelTrans↓ | AbsRot↓ | AbsTrans↓ | Rot† ↓ | RTA15† ↑ | RTA30† ↑ |
|---|---|---|---|---|---|---|---|
| Traj-Crafter [48] | 5.94 | 0.50 | 6.93 | 0.52 | 9.76 | 22.96% | 25.93% |
| ReCamM [2] | 4.26 | **0.32** | 10.08 | 0.34 | 7.49 | 7.61% | 10.20% |
| ReCamM+Aug | 3.66 | 0.43 | 11.74 | 0.46 | 13.88 | 3.89% | 5.93% |
| **SpaceTimePilot (ours)** | **2.71** | 0.33 | **5.63** | 0.34 | **4.09** | **35.19%** | **54.44%** |

† Evaluation based on first-frame camera accuracy.

Table 5. Time-embedding compressor ablation. The proposed time-embedding method, trained with temporal warping on the proposed dataset, yields sharper results overall.

| Time Embedding | PSNR↑ | SSIM↑ | LPIPS↓ |
|---|---|---|---|
| Uniform Sampling | 14.10 | 0.5981 | 0.5039 |
| 1D-Conv | 14.75 | 0.6134 | 0.4878 |
| 1D-Conv + Joint Data | 15.41 | 0.6252 | 0.4830 |
| **1D-Conv +Cam×Time** | **21.16** | **0.7674** | **0.1764** |



Figure 6. **Qualitative comparison of disentangled camera-time control.** In this example, we apply reverse playback (time) and a pan-right camera motion starting from the first-frame pose to a source video (top), whose original camera motion is dolly-in (red to blue). SpaceTimePilot, by explicitly disentangling space and time, achieves correct camera control (red boxes) together with accurate temporal control (green boxes). For TrajectoryCrafter, it first reverses the frames and then apply their method for viewpoint control, resulting in incorrect camera motion. ReCamMaster (with joint-dataset training) is unable to perform temporal control, leading to failure cases.
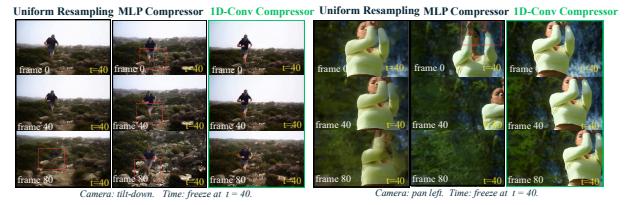


Figure 7. **Temporal compression ablation.** Comparing uniform resampling, MLP, and 1D-Conv compressors under tilt-down and pan-right bullet-time controls, $\mathbf{t}_{trg} = [40, \ldots, 40]$.

resulting in a total of 1800 generated videos. We apply SpatialTracker-v2 [45] to recover camera poses from the generated videos and compare them with the corresponding input camera poses. To ensure consistent scale, we align the magnitude of the first two camera locations. Trajectory accuracy is quantified using **RotErr** and **TransErr** following [8], under two protocols: (1) evaluating the raw trajectories defined w.r.t. the first frame (relative protocol, RelRot, RelTrans) and (2) evaluating after aligning to the estimated pose of the first frame (absolute protocol, AbsRot, AbsTrans). Specifically, we transform the recovered raw trajectories by multiplying the relative pose between the generated and source first frames, estimated by DUSt3R [38]. We also compare this DUSt3R pose with the target trajectory's initial pose, and report RotErr, RTA@15 and RTA@30, as translation magnitude is scale-ambiguous.

To measure only the impact of source camera conditioning, we consider the original ReCamMaster [2] (*ReCamM*) and two variants. Since ReCamMaster is originally trained

on datasets where the first frame of the source and target videos are identical, the model always copies the first frame regardless of the input camera pose. For fairness, we retrain ReCamMaster with more data augmentations to include non-identical first frames, denoted as *ReCamM+Aug*. Next, we condition the model additionally with source cameras $\mathbf{c}_{src}$ following Eq. 3, denoted as *ReCamM+Aug+$\mathbf{c}_{src}$*. Finally we also report the results of TrajectoryCrafter [48].

In Table 4, we observe that the absolute protocol produces consistently higher errors, as trajectories must not only match the overall shape (relative protocol) but also align correctly in position and orientation. Interestingly, ReCamM+Aug yields higher errors than the original ReCamM, whereas incorporating source cameras $\mathbf{c}_{src}$ results in the best overall performance. This suggests that, without explicit reference to $\mathbf{c}_{src}$, exposure to more augmented videos with differing initial frames can instead confuse the model. The newly introduced conditioning signal on the source video's trajectory $\mathbf{c}_{src}$ achieves substantially better camera-control accuracy across all metrics, more reliable first-frame alignment, and more faithful adherence to the full trajectory than all baselines.

#### 4.1.4. Qualitative results.

Besides the quantitative evaluation, we also demonstrate the strength of SpaceTimePilot with visual examples. In Fig. 6, we show that only our method correctly synthesizes both the camera motion (red boxes) and the animation-time state (green boxes). While ReCamMaster handles camera control well, it cannot modify the temporal state, such as enabling

reverse playback. TrajectoryCrafter, in contrast, is confused by the reverse frame shuffle, causing the camera pose of the last source frame (blue boxes) to incorrectly appear in the first frame of the generated video. More visual results can be found in Fig. 5.

## 4.2. Ablation Study

To validate the effectiveness of the proposed Time embedding module, in Table 5, we follow the time-control evaluation set up in Sec. 4.1.1 and compare our 1D convolutional time embedding against several variants and alternatives discussed in Sec. 3.2.1: (1) Uniform-Sampling: sampling the 81-frame embedding uniformly to a 21-frame sequence, which is equivalent to adopting sinusoidal embeddings at the latent frame $f'$ level; (2) 1D-Conv: using 1D convolution layers to compress from $\mathbf{t} \in \mathbb{R}^F$ to $\mathbf{t} \in \mathbb{R}^{F'}$, trained with ReCamMaster and SynCamMaster datasets. (3) 1D-Conv+jointdata: row 2 but including additionally static-scene datasets [18, 53]. (4) 1D-Conv (ours): row 2 but instead including the proposed Cam×Time. We observe that applying a 1D convolution to learn a compact representation by compressing the fine-grained $F$-dim embeddings into a $F'$-dim space performs notably better than directly constructing sinusoidal embeddings at the coarse $f'$ level. Incorporating static-scene datasets yields only limited improvements, likely due to their restricted temporal control patterns. By contrast, using the proposed Cam×Time consistently delivers the largest gains across all three metrics, confirming the effectiveness of our newly introduced datasets. Furthermore, as shown in Fig. 7, we present a visual comparison of bullet-time results using uniform sampling and an MLP instead of the 1D convolution for compressing the temporal control signal. Uniform sampling produces noticeable artifacts, and the MLP compressor causes abrupt camera motion, whereas the 1D convolution effectively locks the animation time and enables smooth camera movement.

## 5. Conclusion

We present SpaceTimePilot, the first video diffusion model to provide fully disentangled spatial and temporal control, enabling 4D space-time exploration from a single monocular video. Our method introduces a new "animation time" representation together with a source-aware camera-control mechanism that leverages both source and target poses. This is supported by the synthetic Cam×Time and a temporal-warping training scheme, which supply dense spatiotemporal supervision. These components allow precise camera and time manipulation, arbitrary initial poses, and flexible multi-round generation. Across extensive experiments, SpaceTimePilot consistently surpasses state-of-the-art baselines, offering significantly improved camera-control accuracy and reliable execution of complex retiming

effects such as reverse playback, slow motion, and bullet-time.

## 6. Acknowledgement

## References

[1] Jianhong Bai, Menghan Xia, Xintao Wang, Ziyang Yuan, Xiao Fu, Zuozhu Liu, Haoji Hu, Pengfei Wan, and Di Zhang. SynCamMaster: Synchronizing multi-camera video generation from diverse viewpoints. *arXiv preprint arXiv:2412.07760*, 2024. 2, 3, 5, 14, 16, 17

[2] Jianhong Bai, Menghan Xia, Xiao Fu, Xintao Wang, Lianrui Mu, Jinwen Cao, Zuozhu Liu, Haoji Hu, Xiang Bai, Pengfei Wan, and Di Zhang. ReCamMaster: Camera-controlled generative rendering from a single video. In *ICCV*, 2025. 2, 3, 4, 5, 6, 7, 8, 14, 16, 17

[3] Omer Bar-Tal, Hila Chefer, Omer Tov, Charles Herrmann, Roni Paiss, Shiran Zada, Ariel Ephrat, Junhwa Hur, Guanghui Liu, Amit Raj, et al. Lumiere: A space-time diffusion model for video generation. In *SIGGRAPH Asia 2024 Conference Papers*, pages 1–11, 2024. 2

[4] 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.

[5] Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video generation models as world simulators. 2024. 2

[6] Chen Gao, Ayush Saraf, Johannes Kopf, and Jia-Bin Huang. Dynamic view synthesis from dynamic monocular video. In *ICCV*, pages 5712–5721, 2021. 2

[7] Klaus Greff, Francois Belletti, Lucas Beyer, Carl Doersch, Yilun Du, Daniel Duckworth, David J Fleet, Dan Gnanapragasam, Florian Golemo, Charles Herrmann, Thomas Kipf, Abhijit Kundu, Dmitry Lagun, Issam Laradji, Hsueh-Ti (Derek) Liu, Henning Meyer, Yishu Miao, Derek Nowrouzezahrai, Cengiz Oztireli, Etienne Pot, Noha Radwan, Daniel Rebain, Sara Sabour, Mehdi S. M. Sajjadi, Matan Sela, Vincent Sitzmann, Austin Stone, Deqing Sun, Suhani Vora, Ziyu Wang, Tianhao Wu, Kwang Moo Yi, Fangcheng Zhong, and Andrea Tagliasacchi. Kubric: a scalable dataset generator. In *CVPR*, 2022. 3

[8] Hao He, Yinghao Xu, Yuwei Guo, Gordon Wetzstein, Bo Dai, Hongsheng Li, and Ceyuan Yang. CameraCtrl: Enabling camera control for text-to-video generation. In *ICLR*, 2025. 8

[9] Hao He, Ceyuan Yang, Shanchuan Lin, Yinghao Xu, Meng Wei, Liangke Gui, Qi Zhao, Gordon Wetzstein, Lu Jiang, and

Hongsheng Li. CameraCtrl II: Dynamic scene exploration via camera-controlled video diffusion models. *arXiv preprint arXiv:2503.10592*, 2025. 2

[10] 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 *CVPR*, 2024. 7

[11] Adobe Systems Inc. Mixamo, 2018. Accessed: 2025-03-07. 12

[12] Hyeonho Jeong, Chun-Hao Paul Huang, Jong Chul Ye, Niloy Mitra, and Duygu Ceylan. Track4Gen: Teaching video diffusion models to track points improves video generation. In *CVPR*, 2025. 2

[13] Hyeonho Jeong, Suhyeon Lee, and Jong Chul Ye. Reangle-A-Video: 4d video generation as video-to-video translation. In *ICCV*, 2025. 3

[14] 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. In *ICLR*, 2025. 2

[15] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Trans. Graph.*, 42(4):139–1, 2023. 2

[16] Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangfeng Xiong, Xin Li, Bo Wu, Jianwei Zhang, et al. Hunyuanvideo: A systematic framework for large video generative models. *arXiv preprint arXiv:2412.03603*, 2024. 2

[17] Jiahui Lei, Yijia Weng, Adam Harley, Leonidas Guibas, and Kostas Daniilidis. MoSca: Dynamic gaussian fusion from casual videos via 4d motion scaffolds. In *CVPR*, 2025. 2

[18] Zhengqi Li, Tali Dekel, Forrester Cole, Richard Tucker, Noah Snavely, Ce Liu, and William T Freeman. Learning the depths of moving people by watching frozen people. In *CVPR*, 2019. 7, 9, 14, 16

[19] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. Neural scene flow fields for space-time view synthesis of dynamic scenes. In *CVPR*, pages 6498–6508, 2021. 2

[20] Zhengqi Li, Qianqian Wang, Forrester Cole, Richard Tucker, and Noah Snavely. Dynibar: Neural dynamic image-based rendering. In *CVPR*, 2023. 2

[21] Hanwen Liang, Yuyang Yin, Dejia Xu, Hanxue Liang, Zhangyang Wang, Konstantinos N Plataniotis, Yao Zhao, and Yunchao Wei. Diffusion4D: Fast spatial-temporal consistent 4d generation via video diffusion models. In *NeurIPS*, 2024. 2

[22] Jinwei Lin. Dynamic NeRF: A review. *arXiv preprint arXiv:2405.08609*, 2024. 2

[23] 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 *CVPR*, pages 22160–22169, 2024. 2, 5

[24] Jiaxin Lu, Chun-Hao Paul Huang, Uttaran Bhattacharya, Qixing Huang, and Yi Zhou. Humoto: A 4d dataset of mocap human object interactions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 10886–10897, 2025. 12

[25] 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 *ECCV*, pages 99–106. ACM New York, NY, USA, 2021. 2

[26] Kepan Nan, Rui Xie, Penghao Zhou, Tiehan Fan, Zhenheng Yang, Zhijie Chen, Xiang Li, Jian Yang, and Ying Tai. OpenVid-1M: A large-scale high-quality dataset for text-to-video generation. *arXiv preprint arXiv:2407.02371*, 2024. 7

[27] 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 *ICCV*, pages 5865–5874, 2021. 2

[28] Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T. Barron, Sofien Bouaziz, Dan B. Goldman, Ricardo Martin-Brualla, and Steven M. Seitz. HyperNeRF: A higher-dimensional representation for topologically varying neural radiance fields. *ACM Transactions on Graphics (TOG)*, 40 (6):238:1–238:12, 2021. 2

[29] Jack Parker-Holder and Shlomi Fruchter. Genie 3: A new frontier for world models. Google DeepMind Blog, 2025. Accessed: ¡insert date you retrieved¿. 2, 3

[30] Adam Polyak, Amit Zohar, Andrew Brown, Andros Tjandra, Animesh Sinha, Ann Lee, Apoorv Vyas, Bowen Shi, Chih-Yao Ma, Ching-Yao Chuang, David Yan, Dhruv Choudhary, Dingkang Wang, Geet Sethi, Guan Pang, Haoyu Ma, Ishan Misra, Ji Hou, Jialiang Wang, Kiran Jagadeesh, Kunpeng Li, Luxin Zhang, Mannat Singh, Mary Williamson, Matt Le, Matthew Yu, Mitesh Kumar Singh, Peizhao Zhang, Peter Vajda, Quentin Duval, Rohit Girdhar, Roshan Sumbaly, Sai Saketh Rambhatla, Sam Tsai, Samaneh Azadi, Samyak Datta, Sanyuan Chen, Sean Bell, Sharadh Ramaswamy, Shelly Sheynin, Siddharth Bhattacharya, Simran Motwani, Tao Xu, Tianhe Li, Tingbo Hou, Wei-Ning Hsu, Xi Yin, Xiaoliang Dai, Yaniv Taigman, Yaqiao Luo, Yen-Cheng Liu, Yi-Chiao Wu, Yue Zhao, Yuval Kirstain, Zecheng He, Zijian He, Albert Pumarola, Ali Thabet, Artsiom Sanakoyeu, Arun Mallya, Baishan Guo, Boris Araya, Breena Kerr, Carleigh Wood, Ce Liu, Cen Peng, Dimitry Vengertsev, Edgar Schonfeld, Elliot Blanchard, Felix Juefei-Xu, Fraylie Nord, Jeff Liang, John Hoffman, Jonas Kohler, Kaolin Fire, Karthik Sivakumar, Lawrence Chen, Licheng Yu, Luya Gao, Markos Georgopoulos, Rashel Moritz, Sara K. Sampson, Shikai Li, Simone Parmeggiani, Steve Fine, Tara Fowler, Vladan Petrovic, and Yuming Du. Movie gen: A cast of media foundation models, 2024. 2

[31] Xuanchi Ren, Tianchang Shen, Jiahui Huang, Huan Ling, Yifan Lu, Merlin Nimier-David, Thomas Müller, Alexander Keller, Sanja Fidler, and Jun Gao. GEN3C: 3d-informed world-consistent video generation with precise camera control. In *CVPR*, 2025. 3

[32] Chris Rockwell, Joseph Tung, Tsung-Yi Lin, Ming-Yu Liu, David F. Fouhey, and Chen-Hsuan Lin. Dynamic camera poses and where to find them. In *CVPR*, 2025. 5

[33] Basile Van Hoorick, Rundi Wu, Ege Ozguroglu, Kyle Sargent, Ruoshi Liu, Pavel Tokmakov, Achal Dave, Changxi Zheng, and Carl Vondrick. Generative camera dolly: Extreme monocular dynamic novel view synthesis. *ECCV*, 2024. 2, 3, 4, 5

[34] Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu, Haiming Zhao, Jianxiao Yang, Jianyuan Zeng, Jiayu Wang, Jingfeng Zhang, Jingren Zhou, Jinkai Wang, Jixuan Chen, Kai Zhu, Kang Zhao, Keyu Yan, Lianghua Huang, Mengyang Feng, Ningyi Zhang, Pandeng Li, Pingyu Wu, Ruihang Chu, Ruili Feng, Shiwei Zhang, Siyang Sun, Tao Fang, Tianxing Wang, Tianyi Gui, Tingyu Weng, Tong Shen, Wei Lin, Wei Wang, Wei Wang, Wenmeng Zhou, Wente Wang, Wenting Shen, Wenyuan Yu, Xianzhong Shi, Xiaoming Huang, Xin Xu, Yan Kou, Yangyu Lv, Yifei Li, Yijing Liu, Yiming Wang, Yingya Zhang, Yitong Huang, Yong Li, You Wu, Yu Liu, Yulin Pan, Yun Zheng, Yuntao Hong, Yupeng Shi, Yutong Feng, Zeyinzi Jiang, Zhen Han, Zhi-Fan Wu, and Ziyu Liu. Wan: Open and advanced large-scale video generative models. *arXiv preprint arXiv:2503.20314*, 2025. 4, 6

[35] Chaoyang Wang, Ashkan Mirzaei, Vidit Goel, Willi Menapace, Aliaksandr Siarohin, Avalon Vinella, Michael Vasilkovsky, Ivan Skorokhodov, Vladislav Shakhrai, Sergey Korolev, Sergey Tulyakov, and Peter Wonka. 4real-video-v2: Fused view-time attention and feedforward reconstruction for 4d scene generation. In *Adv. Neural Inform. Process. Syst.*, 2025. 3

[36] Chaoyang Wang, Peiye Zhuang, Tuan Duc Ngo, Willi Menapace, Aliaksandr Siarohin, Michael Vasilkovsky, Ivan Skorokhodov, Sergey Tulyakov, Peter Wonka, and Hsin-Ying Lee. 4real-video: Learning generalizable photo-realistic 4d video diffusion. In *CVPR*, 2025. 3

[37] Qianqian Wang, Vickie Ye, Hang Gao, Jake Austin, Zhengqi Li, and Angjoo Kanazawa. Shape of motion: 4d reconstruction from a single video. In *ICCV*, 2025. 2

[38] Shuzhe Wang, Vincent Leroy, Yohann Cabon, Boris Chidlovskii, and Jerome Revaud. DUSt3R: Geometric 3d vision made easy. In *CVPR*, 2024. 8

[39] Yaohui Wang, Xinyuan Chen, Xin Ma, Shangchen Zhou, Ziqi Huang, Yi Wang, Ceyuan Yang, Yinan He, Jiashuo Yu, Peiqing Yang, et al. Lavie: High-quality video generation with cascaded latent diffusion models. *IJCV*, pages 1–20, 2024. 2

[40] Zun Wang, Jaemin Cho, Jialu Li, Han Lin, Jaehong Yoon, Yue Zhang, and Mohit Bansal. EPiC: Efficient Video Camera Control Learning with Precise Anchor-Video. *arXiv preprint arXiv:2505.21876*, 2025. 3

[41] Daniel Watson, Saurabh Saxena, Lala Li, Andrea Tagliasacchi, and David J. Fleet. Controlling space and time with diffusion models. In *ICLR*, 2025. 2, 3, 4, 5, 7

[42] Guanjun Wu, Taoran Yi, Jiemin Fang, Lingxi Xie, Xiaopeng Zhang, Wei Wei, Wenyu Liu, Qi Tian, and Xinggang Wang. 4d gaussian splatting for real-time dynamic scene rendering.

In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 20310–20320, 2024. 2

[43] Rundi Wu, Ruiqi Gao, Ben Poole, Alex Trevithick, Changxi Zheng, Jonathan T Barron, and Aleksander Holynski. Cat4D: Create anything in 4d with multi-view video diffusion models. In *CVPR*, 2024. 2, 3, 4, 5, 7

[44] Tong Wu, Shuai Yang, Ryan Po, Yinghao Xu, Ziwei Liu, Dahua Lin, and Gordon Wetzstein. Video world models with long-term spatial memory. *arXiv preprint arXiv:2506.05284*, 2025. 3

[45] Yuxi Xiao, Jianyuan Wang, Nan Xue, Nikita Karaev, Iurii Makarov, Bingyi Kang, Xin Zhu, Hujun Bao, Yujun Shen, and Xiaowei Zhou. SpatialTrackerV2: 3d point tracking made easy. In *ICCV*, 2025. 8

[46] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion models with an expert transformer. *arXiv preprint arXiv:2408.06072*, 2024. 2

[47] Jae Shin Yoon, Kihwan Kim, Orazio Gallo, Hyun Soo Park, and Jan Kautz. Novel view synthesis of dynamic scenes with globally coherent depths from a monocular camera. In *CVPR*, pages 5339–5348, 2020. 2

[48] Mark YU, Wenbo Hu, Jinbo Xing, and Ying Shan. TrajectoryCrafter: Redirecting camera trajectory for monocular videos via diffusion models. In *ICCV*, 2025. 3, 7, 8

[49] 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. 3

[50] David Junhao Zhang, Roni Paiss, Shiran Zada, Nikhil Karnad, David E Jacobs, Yael Pritch, Inbar Mosseri, Mike Zheng Shou, Neal Wadhwa, and Nataniel Ruiz. ReCapture: Generative video camera controls for user-provided videos using masked video fine-tuning. In *CVPR*, 2024. 3

[51] David Junhao Zhang, Jay Zhangjie Wu, Jia-Wei Liu, Rui Zhao, Lingmin Ran, Yuchao Gu, Difei Gao, and Mike Zheng Shou. Show-1: Marrying pixel and latent diffusion models for text-to-video generation. *IJCV*, pages 1–15, 2024. 2

[52] Jensen (Jinghao) Zhou, Hang Gao, Vikram Voleti, Aaryaman Vasishta, Chun-Han Yao, Mark Boss, Philip Torr, Christian Rupprecht, and Varun Jampani. Stable Virtual Camera: Generative view synthesis with diffusion models. *arXiv preprint*, 2025. 3

[53] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: Learning view synthesis using multiplane images. In *SIGGRAPH*, 2018. 2, 5, 7, 9, 14, 16

## A. Network Architecture

The network architecture of SpaceTimePilot is depicted in Fig. 8. The newly introduced animation-time embedder $\mathcal{E}_{\text{ani}}$ encodes the source and target animation times, $\mathbf{t}_{\text{src}}$ and $\mathbf{t}_{\text{trg}}$, into tensors matching the shapes of $x_{\text{src}}$ and $x_{\text{trg}}$, which are then added to them respectively. During training, we train only the camera embedder $\mathcal{E}_{\text{cam}}$, the animation-time embedder $\mathcal{E}_{\text{ani}}$, the self-attention (full-3D attention), and the projector layers before the cross-attention.



Figure 8. **Architecture of SpaceTimePilot.** Our model jointly conditions on camera trajectories and temporal control signals via space–time attention, enabling non-monotonic motion generation such as reversals, repeats, accelerations, and zigzag time.

## B. Longer Space-Time Exploration Video with Disentangled Controls

One of the central advantages of SpaceTimePilot is its ability to freely navigate both spatial and temporal dimensions, with arbitrary starting points in each dimension and fully customizable trajectories through them. Although each individual generation is limited to an 81-frame window, we show that SpaceTimePilot can effectively extend this window indefinitely through a multi-turn autoregressive inference scheme, enabling continuous and controllable space–time exploration from a single input video. The overall pipeline is illustrated in Fig. 9.

The core idea is to generate the final video in autoregressive segments that connect seamlessly. For example, given a source video of 81 frames, we may first generate a $0.5\times$ slow-motion sequence covering frames 0–40 with a new camera trajectory. Then, continuing both the visual context and the generated camera trajectory, we can produce the next segment starting from the final camera pose

of the previous output, while temporally traversing the remaining frames 40–81. This yields an autoregressive chain of viewpoint-controlled video segments that together create a continuous long-range space–time trajectory.

A key property that enables this behavior is that our model can generate video segments whose camera poses do *not* need to start at the first frame. This allows precise control over the starting point, both in time and viewpoint, for every generated chunk, ensuring smooth, consistent motion over extended sequences.

To maintain contextual coherence across iterations, we introduce a lightweight memory mechanism. During training, the model is conditioned on **a pair of source videos**, which enables consistent chaining during inference. Specifically:

- At iteration $i = 1$, the model is conditioned only on the original source video.
- At iteration $i = 2$, it is conditioned on both the source video and the previously generated 81-frame segment.
- This process repeats, with each iteration conditioning on the source video as well as the most recent generated segment.

This simple yet effective strategy allows SpaceTimePilot to generate arbitrarily long, smoothly connected sequences with continuous and precise control over both temporal manipulation and camera motion.

Here, we showcase how this can be used to conduct large viewpoint changes, as demonstrated in Fig. 10.

## C. Additional Details on the Proposed Cam×Time Dataset.

The `Cam×Time` dataset is built using high-quality, commercially licensed 3D environments that include both realistic indoor and outdoor scenes. For each environment, we populate the space with multiple animated human characters. The character assets are sourced from Mixamo [11] and HUMOTO [24], and each character is manually textured and refined to ensure realistic geometry, appearance, and material quality. The animations span a diverse range of human motions, including locomotion, gestures, and human-object interactions. Examples of scenes are shown in Fig. 11. Please refer to the complementary website for the video examples.

To capture rich spatial coverage, we generate four distinct camera trajectories for every scene. Camera paths include rotational orbits, linear tracking motions, and smoothly curved arcs. A dedicated validity module ensures that each trajectory: (1) begins at a collision-free location with clear visibility of the main character, (2) maintains non-intersecting movement with the environment throughout the path, and (3) preserves full subject visibility across all viewpoints.
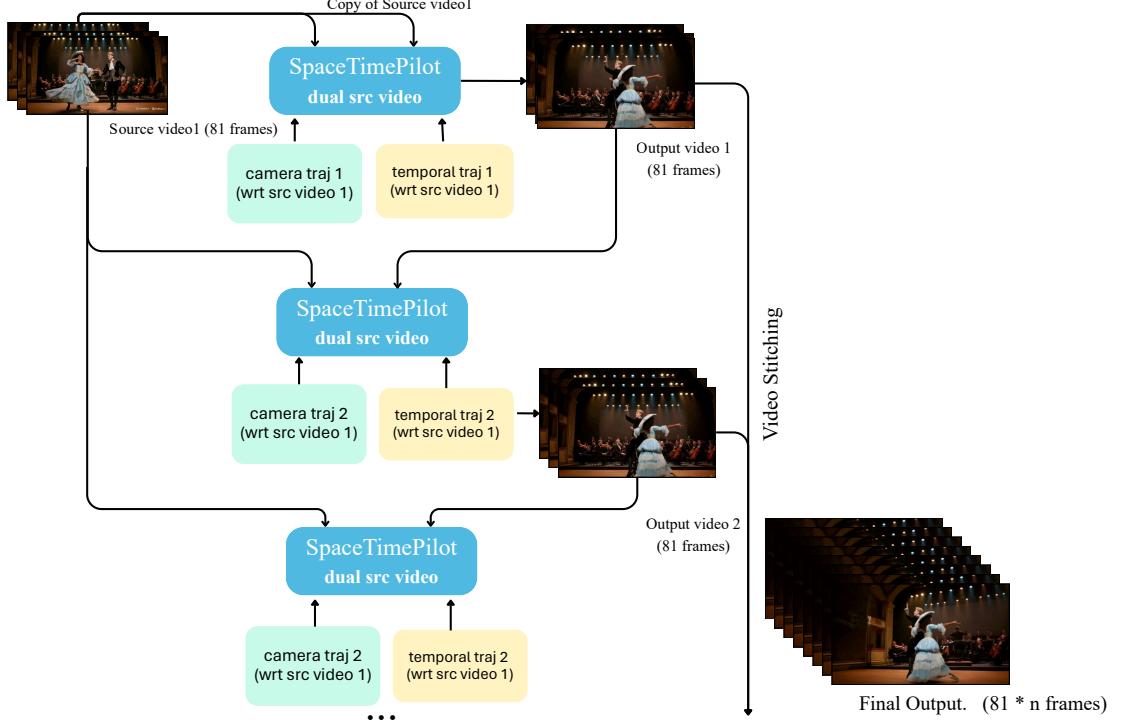
Figure 9. **Overview of the multi-turn autoregressive inference scheme.** The model first generates an 81-frame segment conditioned on the source video and a chosen space–time trajectory. The resulting output is then reused as a secondary source video for subsequent iterations, each with its own camera and temporal trajectory. By chaining these iterations and stitching the outputs, SpaceTimePilot produces a long, coherent video that follows an arbitrary space–time path.



Figure 10. **Multi-turn autoregressive generation with SpaceTimePilot.** Top row: source video frames. Rows 2–4: Turn-1, Turn-2, and Turn-3 generations. At each turn, SpaceTimePilot jointly conditions on (1) the original source video and (2) the previously generated chunk, ensuring temporal continuity, stable motion progression, and consistent camera geometry. This dual-conditioning design enables viewpoint changes far beyond the input video—such as rotating to the rear of the tiger or transitioning from a low-angle shot to a high bird's-eye view—while preserving visual and motion coherence. Please refer to section "AR Demos" in the website for videos.

Each trajectory is rendered into a 120-frame sequence at a resolution of $1080 \times 1080$ pixels, providing dense temporal sampling with high visual fidelity. This yields three multi-view video sequences per scene, each covering the full motion duration with consistent lighting, textures, and geometry. Overall, we rendered 1500 videos from 500 animations, each with 120 videos full grid rendering, leading to 180k videos.

For temporal-control training, we could sample any time variants from these sequences, including slow motion, reverse playback, bullet-time around arbitrary frames, and non-monotonic time patterns such as forward–backward oscillation. These augmented temporal signals are illustrated in Fig. 12.

Figure 11. **Example of Cam×Time**. Multi-view, densely sampled sequences from the Cam×Time dataset. Each row shows frames from one camera trajectory, and each column samples different timesteps (0–120). The dataset provides diverse environments, human motions, and four camera paths per scene with full 120-frame temporal coverage.
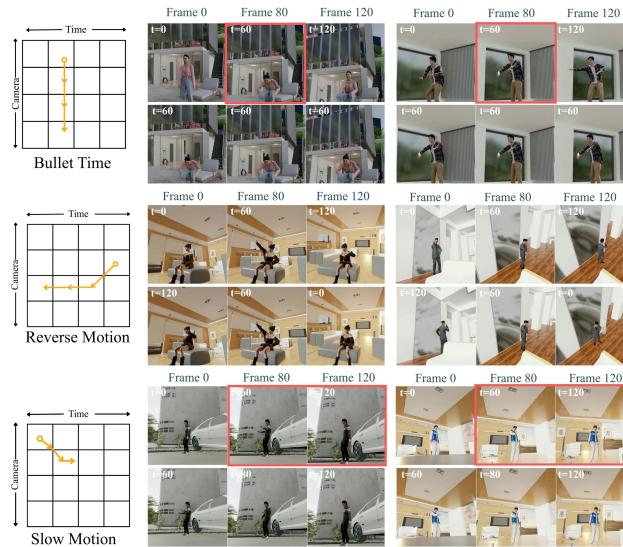


Figure 12. **Sampling from Cam×Time.** By sampling from the Cam×Time dataset, we can extract frames corresponding to arbitrary combinations of camera viewpoints and temporal positions, forming source-target pairs with rich camera and temporal control signals.

# D. Additional Ablation Studies

## D.1. Temporal Warping Augmentation

Using [1, 2] as our default datasets, we compare training jointly with static-scene datasets [18, 53] with applying only temporal warping (TW) augmentation on the default datasets (Sec. 3.2.2 in the main paper). Although static-

scene datasets naturally support bullet-time effects, they do not provide enough diversity of temporal control configurations for models to reliably learn time locking on their own, as shown in Fig. 14 (top). Please refer to section "Effective Temporal Warping" in the website for more videos.

In Fig. 14 (bottom), we further show that freezing temporal warping (3rd row) produces better results than training without freezing it. Please refer to section "Freeze Warping Ablations" in the website for more videos.

## D.2. Significance of Cam×Time **Dataset**

Besides the quantitative results in the main paper (Table 5), in Fig. 15 (top), we provide visual comparisons demonstrating the effectiveness of the proposed Cam×Time dataset. Clear artifacts appear in baselines trained without additional data or with only static-scene augmentation (highlighted in red boxes), whereas incorporating Cam×Time removes these artifacts, demonstrating its significance. Please refer to section "Dataset Ablations" in the website for more videos.

## D.3. Time Embedding Ablation

As promised in Sec. 3.2.1 in the main paper, we compare several time-embedding strategies. RoPE($f'$) can freeze the scene dynamics at $t$=40, but it also undesirably locks the camera motion. Using MLP, by contrast, fails to lock the temporal state at all (red boxes). Conditioning on the latent frame $f'$ (with uniform sampling) introduces noticeable artifacts. In comparison, the proposed 1D-Conv embedding enables SpaceTimePilot to preserve the intended scene dynamics while still generating accurate camera motion. Adding Cam×Time to training further enhances the results. Please refer to section "Time-Embedding Method Ablation" in the website for more examples.

# E. Additional Qualitative Visualizations

We show more qualitative results of SpaceTimePilot in Fig. 13. Our model provides fully disentangled control over camera motion and temporal dynamics. Each row presents a different pairing of temporal control inputs (top-left icon) and camera trajectories. SpaceTimePilot reliably generates coherent videos under diverse conditions, including normal and reverse playback, bullet-time, slow motion, replay motion, and complex camera movements such as panning, tilting, zooming, and vertical translation. Please refer to section "Video Demonstrations" in the website for more examples.
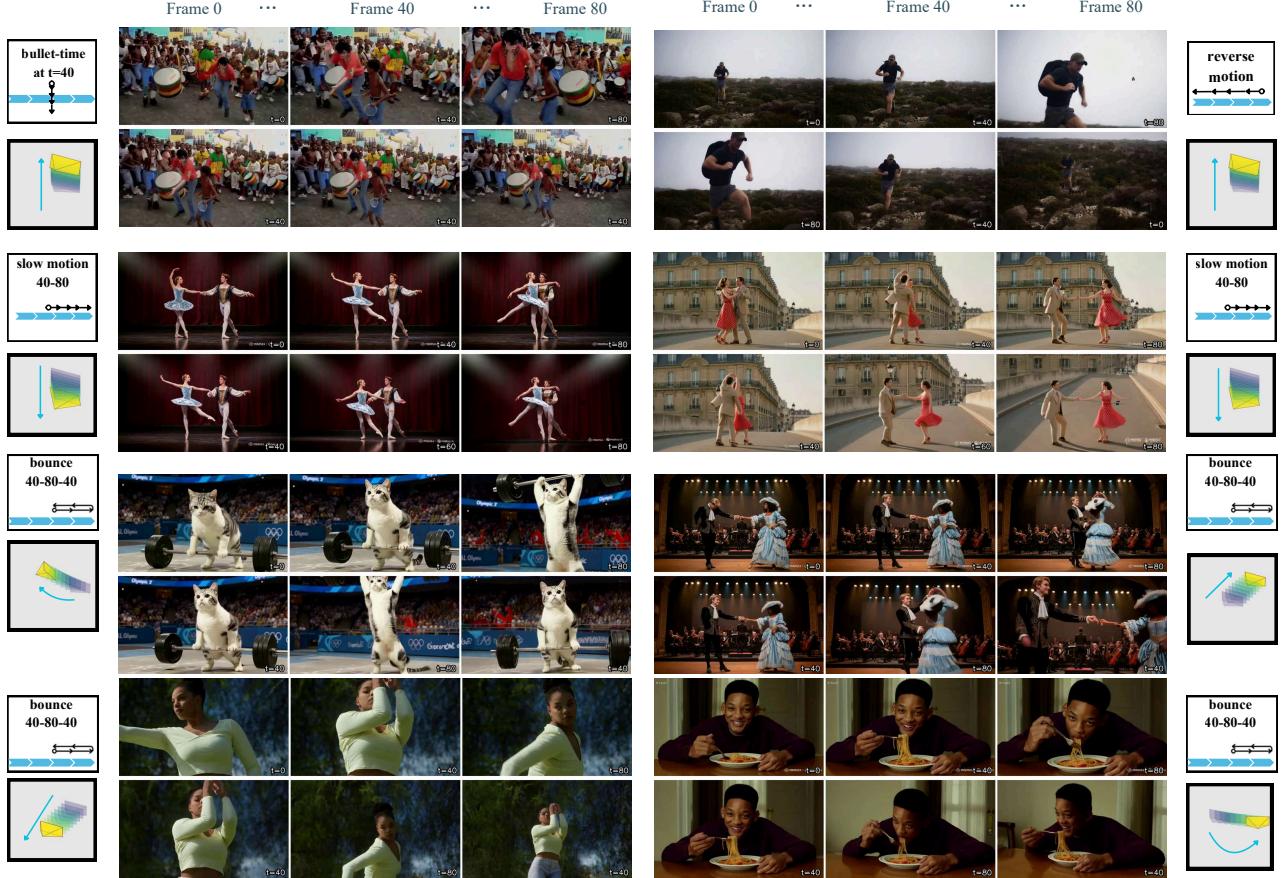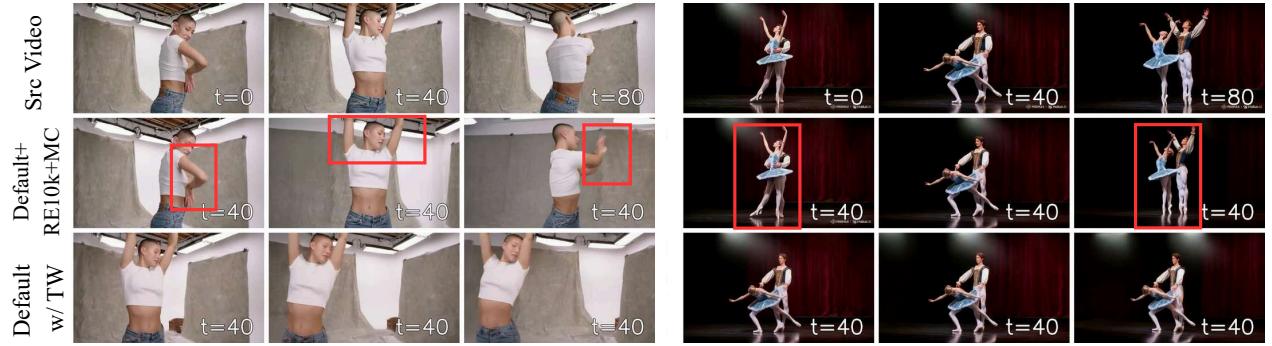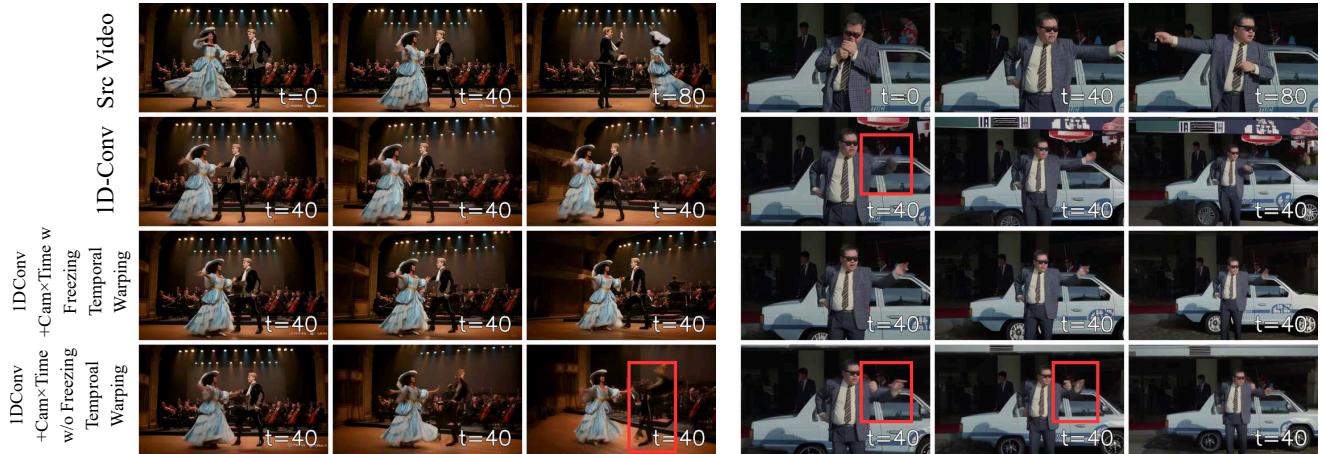
Figure 13. **More Qualitative results.** Our model provides fully disentangled control over camera motion and temporal dynamics. Each row illustrates a different combination of temporal control inputs (top-left icon) and camera trajectories. SpaceTimePilot consistently produces coherent videos across a wide range of controls, including normal and reverse playback, bullet-time, slow motion, replay motion, and complex camera paths such as panning, tilting, zooming, and vertical motion.

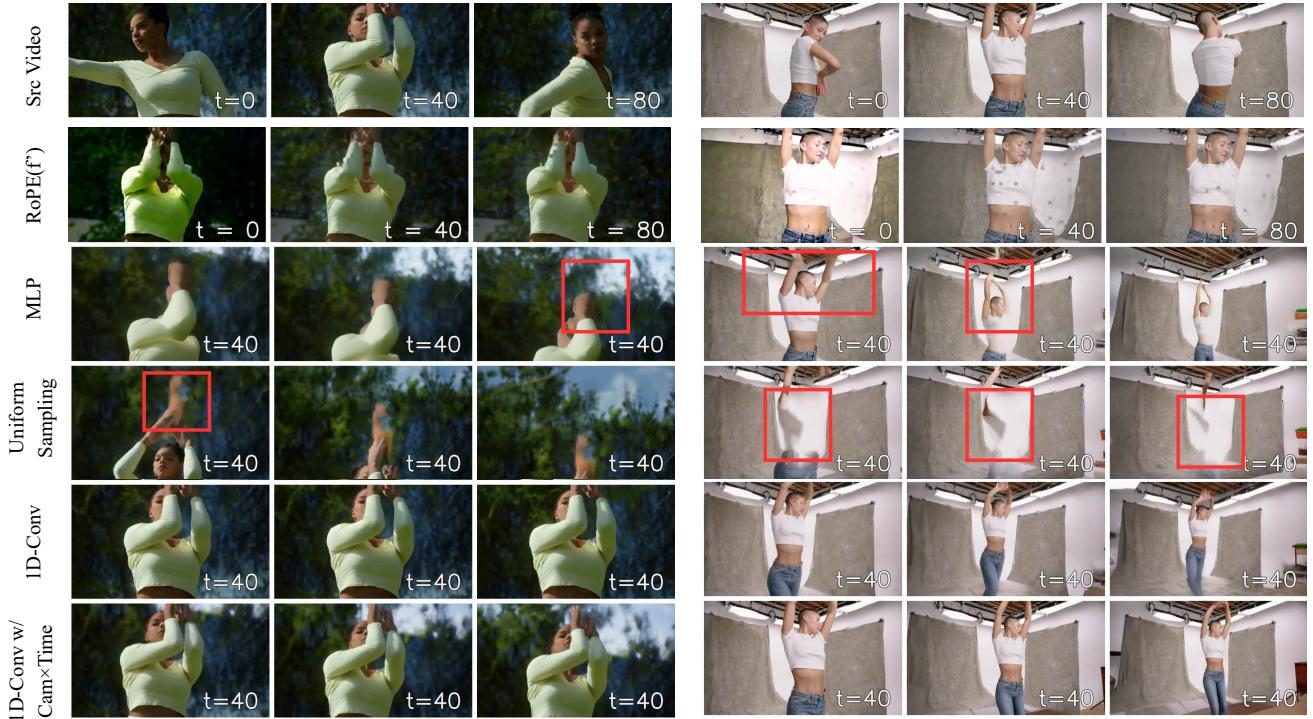Ablation 1: Temporal warping (TW) vs. Joint dataset training on bullet-time effect (**t**=40)



Ablation 2: Varied temporal warping configuration on bullet-time effect (**t**=40)

Figure 14. **Ablation study.** (Top) Using [1, 2] as default datasets, we compare the influence of adding static-scene datasets [18, 53] vs. just doing temporal warping (TW) augmentation (Sec. 3.2.2 in the main paper). Temporal warping definitely provide more variety of time control signals, allowing models to learn better camera-time disentanglement. (Bottom) We further compare different configurations of warping, where we show freezing temporal warping (3$^{\text{rd}}$ row) leads to better results than those trained without freezing temporal warping.

Ablation 3: Dataset ablation: Cam×Time vs. static-scene datasets on bullet-time effect (**t**=40)



Ablation 4: Different time embedding schemes on bullet-time effect (**t**=40)

Figure 15. **Ablation study.** (Top) We verify the efficacy of the proposed Cam×Time dataset. Considering [1, 2] as default datasets, we compare the impact of different datasets on the generated videos. One can clearly see artifacts in baselines without any extra data or augmented with static-scene data, whereas training additionally with Cam×Time leads to no artifacts, confirming the usefulness of our dataset. (Bottom) We compare several time-embedding strategies. The MLP fails to lock the temporal state (red boxes), while RoPE($f'$) correctly freezes the scene dynamics at **t**=40 but unintentionally locks the camera motion too. Conditioning on the latent frame $f'$ (with uniform sampling) introduces noticeable artifacts. In contrast, the proposed 1D-Conv embedding allows SpaceTimePilot to both freeze the scene dynamics at **t**=40 and produce intended camera motion. Incorporating Cam×Time during training further improves performance.