# Infinite-Homography as Robust Conditioning for Camera-Controlled Video Generation

Min-Jung Kim*    Jeongho Kim*    Hoiyeong Jin    Junha Hyung    Jaegul Choo
KAIST AI

{emjay73,rlawjdghek,hy.jin,sharpeeee,jchoo}@kaist.ac.kr
*Equal contribution

Figure 1. **InfCam Results.** Given a video and a target camera trajectory, InfCam generates a video that faithfully follows the specified camera path. The world coordinate origin is defined by the first frame's camera pose (highlighted in red). The leftmost column visualizes the backward, arc, and rotational camera trajectories, and the right side shows input–generated video pairs corresponding to each trajectory. The rotational trajectory is generated with a shorter focal length to illustrate wide field-of-view generation. The black dashed box in the last row indicates the original field-of-view of the input video.

## Abstract

*Recent progress in video diffusion models has spurred growing interest in camera-controlled novel-view video generation for dynamic scenes, aiming to provide creators with cinematic camera control capabilities in post-production. A key challenge in camera-controlled video generation is ensuring fidelity to the specified camera pose, while maintaining view consistency and reasoning about occluded geometry from limited observations. To address this, existing methods either train trajectory-conditioned video generation model on trajectory–video pair dataset, or estimate depth from the input video to reproject it along a target trajectory and generate the unprojected regions. Nevertheless, existing methods struggle to generate camera-pose–faithful, high-quality videos for two main reasons: (1) reprojection-based approaches are highly susceptible to errors caused by inaccurate depth estimation; and (2) the*

*limited diversity of camera trajectories in existing datasets restricts learned models. To address these limitations, we present InfCam, a depth-free, camera-controlled video-to-video generation framework with high pose fidelity. The framework integrates two key components: (1) infinite homography warping, which encodes 3D camera rotations directly within the 2D latent space of a video diffusion model. Conditioning on this noise-free rotational information, the residual parallax term is predicted through end-to-end training to achieve high camera-pose fidelity; and (2) a data augmentation pipeline that transforms existing synthetic multiview datasets into sequences with diverse trajectories and focal lengths. Experimental results demonstrate that InfCam outperforms baseline methods in camera-pose accuracy and visual fidelity, generalizing well from synthetic to real-world data. Link to our project page:* https://emjay73.github.io/InfCam/

## 1. Introduction

Changing the camera viewpoint in post-production is a highly sought-after video editing technique, as it eliminates the need for costly reshoots and enables creative visual effects. Recent advancements in large-scale video diffusion models have spurred research on controlling camera trajectories for novel-view video generation [1, 2, 13, 16, 20].

However, existing methods still suffer from several limitations. First, approaches that explicitly incorporate depth projection [13, 20, 21] suffer from inherent performance degradation due to reprojection errors arising from inaccurately estimated depths. While reprojection serves as an effective mechanism for aligning generated views with the target trajectory when accurate depth maps are available, it can also become a major source of error when the predicted depths are inaccurate. Consequently, the overall performance of such systems remains fundamentally constrained by the accuracy and reliability of the underlying depth predictor. Second, methods that learn the relationship between camera trajectory variations and novel-view video generation from trajectory–video pair datasets [1, 2, 16] tend to inherit dataset-specific biases. When the training data contain biased input–output video pairs, the network is likely to internalize these biases, leading to degraded generalization performance.

To address these limitations, we propose a robust infinite homography-based conditioning method. Unlike existing depth reprojection-based methods that suffer from cascading errors due to their reliance on noisy depth estimates, our infinite homography warping module decomposes the reprojection process into known rotation and unknown translation components. We condition the model on the noise-free rotation component, while the residual translation component is learned end-to-end by leveraging strong

prior knowledge from large video generation models. Additionally, we introduce a data augmentation strategy that transforms datasets with constrained viewpoints into flexible trajectory formats while simultaneously augmenting focal lengths. By synthesizing diverse camera movements and field-of-view variations, our model learns robust view synthesis. Extensive experiments demonstrate that the proposed infinite homography warping module and data augmentation enable superior trajectory fidelity compared to state-of-the-art methods. Fig. 1 demonstrates the generalizability of our approach on challenging real-world videos. Our contributions are summarized as follows:

- We propose a camera-controlled novel-view video generation framework that achieves high pose fidelity.
- We introduce an infinite homography warping module that conditions on a noise-free rotation term, predicting the residual parallax term through end-to-end training.
- We present a data augmentation strategy that transforms existing datasets with limited viewpoints into diverse training pairs with flexible camera trajectories and varying focal lengths, enabling robust generalization
- Extensive experiments demonstrate that our infinite homography warping and data augmentation strategy significantly improve camera-controlled video generation, achieving state-of-the-art performance on both synthetic and real-world videos.

## 2. Related Work

**Depth Reprojection-Based Video Generation.** This line of work [11, 13, 20, 21] estimates per-frame depth from an input and reprojects the resulting 3D structure along a target camera trajectory. The newly exposed or unprojected regions are subsequently predicted or inpainted to produce the final rendered video. For instance, GEN3C [13] proposes a method that builds a spatio-temporal 3D cache from videos via depth estimation and unprojection, renders this cache along the user-specified camera trajectory, and conditions a video diffusion model on these renderings to achieve precise camera control and temporal consistency. Meanwhile, TrajectoryCrafter [20] introduces a dual-stream conditional video diffusion framework that takes as input both the source video and reprojection results, which are subsequently processed by a video inpainting model fine-tuned on a dataset generated using a double reprojection scheme. Since this line of work relies on external depth estimation, its performance is ultimately limited by the accuracy of the estimated depth, which directly affects the magnitude of reprojection errors.

**Trajectory-Conditioned Video Generation.** This line of work [1, 2, 9, 14, 16, 17, 19] trains a video generation model that is explicitly conditioned on camera trajectories, using paired datasets of trajectories and corresponding videos. SynCamMaster [2] builds upon a pre-trained text-
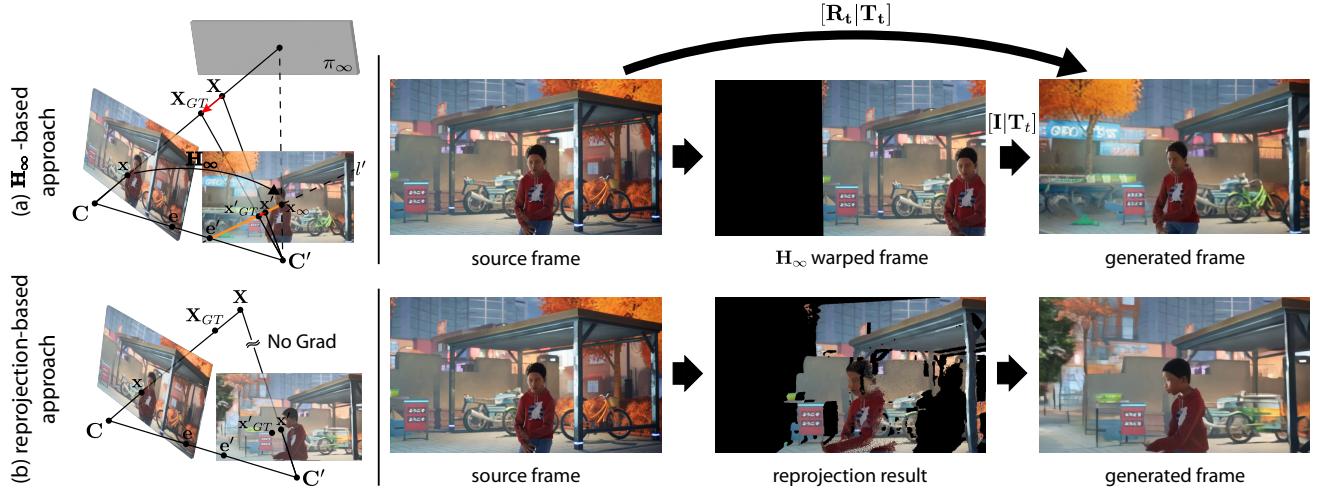
Figure 2. (a) **Infinite homography-based approach (ours).** By conditioning on images warped by $\mathbf{H}\infty$, the model focuses on learning the parallax relative to the plane at infinity. This parallax is confined to the region between the epipole $\mathbf{e}'$ and $\mathbf{x}_\infty$, as visualized by the yellow segment on the epipolar line $l'$. This spatial constraint helps the model to achieve higher camera pose fidelity with reduced search space. End-to-end training enables the network to implicitly refine the 3D geometry, correcting inaccuracies in $\mathbf{X}$. (b) **Reprojection-based approach.** Inaccuracies in depth estimation lead to unreliable conditions, causing artifacts in the generated image. Since no gradients flow through the depth estimation network, the incorrect reprojection position $\mathbf{x}'$ remains fixed during training, hindering error correction.

to-video model with a multi-view synchronization module to produce synchronized videos from diverse camera viewpoints. Trained with a hybrid data scheme and the released SynCamVideo dataset, it demonstrates open-domain novel-view video generation for stationary cameras. ReCamMaster [1] further leverages the generative capabilities of pre-trained text-to-video models through frame-dimension conditioning, and is trained with the MultiCamVideo dataset, a large-scale multi-camera synchronized video dataset built with Unreal Engine 5. While effective, such trajectory-conditioned video generation methods trained on trajectory–video pair data often inherit biases from the underlying training data distribution.

## 3. Motivation

Fig. 2 illustrates the difference between the infinite homography-based approach and the reprojection-based approach at the frame level. For a pair of cameras with centers $\mathbf{C}$ and $\mathbf{C}'$, the image from $\mathbf{C}$ is used to generate the corresponding view of camera $\mathbf{C}'$.

**Reprojection-based Approach.** As illustrated in Fig. 2 (b), The reprojection-based approach first estimates the depth for a pixel $\mathbf{x}$ in camera $\mathbf{C}$, then reproject it to camera $\mathbf{C}'$, producing a corresponding position $\mathbf{x}'$ and yielding a reprojection image. However, any inaccuracy in the predicted depth results in an unreliable condition being provided to the generative model, leading to artifacts in the synthesized image. Moreover, because no gradients flow through the depth estimation network, the inaccurate repro-

jection position $\mathbf{x}'$ remains fixed as a condition throughout training. Although the model can be trained to compensate for these errors, the challenge lies in handling noisy conditions, where it is unclear which parts are corrupted and to what extent, making the problem difficult to address.

**Infinite Homography-based Approach.** The infinite homography $\mathbf{H}_\infty$ represents the homography induced by the plane at infinity $\pi_\infty$ [5]. Given source and target camera intrinsic matrices $\mathbf{K}_s$ and $\mathbf{K}_t$, rotation matrix $\mathbf{R}$, translation vector $\mathbf{t}$, and normal $\mathbf{n}$ of a plane, the infinite homography $\mathbf{H}_\infty$ can be derived from the plane-induced homography $\mathbf{H} = \mathbf{K}_t(\mathbf{R} - \mathbf{t}\mathbf{n}^T/d)\mathbf{K}_s^{-1}$ by taking the limit as the distance $d$ to the plane approaches infinity:

$$\mathbf{H}_\infty = \lim_{d \to \infty} \mathbf{H} = \mathbf{K}_t\mathbf{R}\mathbf{K}_s^{-1}. \quad (1)$$

For a pixel $\mathbf{x}$ with known depth $Z$ measured from the source camera, the reprojection $\mathbf{x}'$ to the target image is expressed as:

$$\mathbf{x}' = \mathbf{K}_t\mathbf{R}\mathbf{K}_s^{-1}\mathbf{x} + \mathbf{K}_t\mathbf{t}/Z = \mathbf{H}_\infty\mathbf{x} + \mathbf{K}_t\mathbf{t}/Z \quad (2)$$

Notably, $\mathbf{H}_\infty$ does not depend on the translation or depth, enabling correspondence of image points at arbitrary depths when the camera undergoes pure rotation ($\mathbf{t} = \mathbf{0}$). The term $\mathbf{K}_t\mathbf{t}/Z$ represents the parallax relative to the plane at infinity. As shown in Fig. 2(a), this parallax is restricted to the region between the epipole $\mathbf{e}'$ and $\mathbf{x}_\infty$, visualized by the yellow segment on the epipolar line $l'$.

Assuming that $\mathbf{K}_s$, $\mathbf{K}_t$, $\mathbf{R}$, and $\mathbf{t}$ are known, conditioning the network on the image warped by $\mathbf{H}_\infty$ allows the
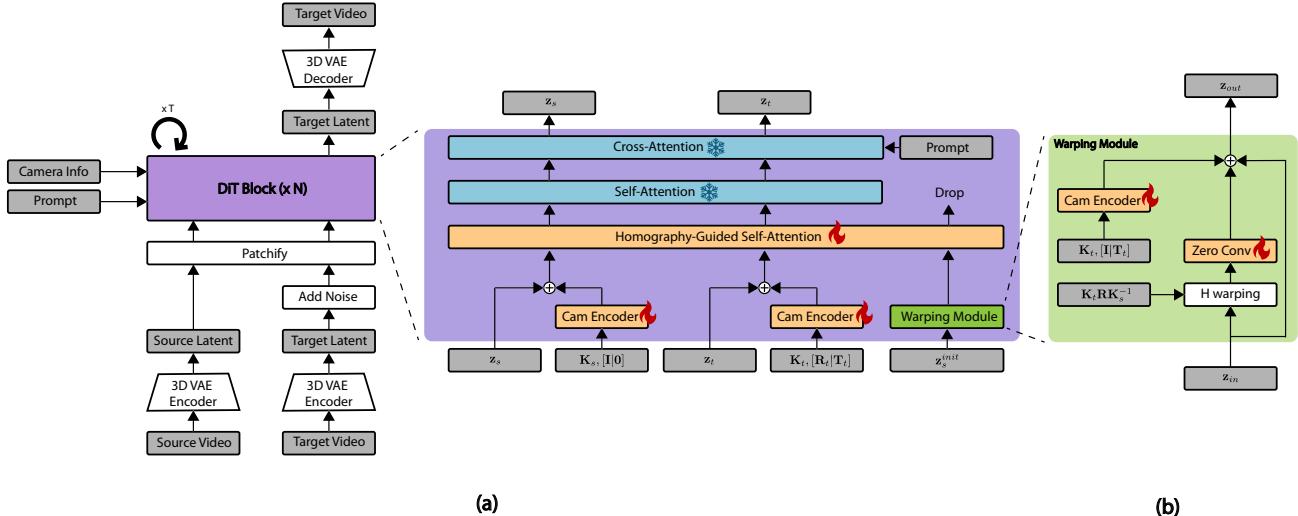
Figure 3. **Model Architecture Overview.** Our model builds upon Wan2.1, training only newly introduced parameters while freezing pretrained weights. (a) **DiT block with homography-guided self-attention layer.** Homography-guided self-attention layer takes source, target, and warped latents combined with camera embeddings as input, and performs per-frame attention, ensuring temporal alignment. By conditioning on warped latents, the model enables rotation-aware reasoning and constrained parallax estimation. Only source and target latents proceed to the subsequent Wan2.1 layers. (b) **Warping module**. This module warps the input latent with infinite homography to handle rotation, then add camera embeddings for translation. This decomposition simplifies reprojection to parallax estimation relative to plane at infinity, enabling higher camera trajectory fidelity.

network to refer to a noise-free image under rotational transformation, enabling it to focus on learning residual information, including $\mathbf{K}_t \mathbf{t}$ and the implicitly predicted depth $Z$, which modulates the final novel views. Since parallax variations are constrained between the epipole $\mathbf{e}'$ and $\mathbf{x}_\infty$, this provides the model with strict boundary conditions (Fig. 2 (a) orange line). This spatial constraint helps the model to improve camera-pose fidelity by effectively narrowing the search space. Furthermore, unlike the reprojection-based approach, our method is trained in an end-to-end manner, enabling the network to implicitly refine the predicted depth $Z$ toward the ground truth as it learns to produce accurate parallax. For computational efficiency, we warp the latent feature rather than the image itself during training.

## 4. Method

Our goal is to perform novel-view video generation that faithfully follows a given target camera trajectory. Specifically, given a source video $\mathbf{V}_s \in \mathbb{R}^{F \times C \times H \times W}$, source and target camera intrinsics $\mathbf{K}_s, \mathbf{K}_t \in \mathbb{R}^{3 \times 3}$, and target camera trajectory $\mathbf{T} \in \mathbb{R}^{F \times 3 \times 4}$ expressed relative to the source video's initial camera pose, our InfCam generates novel view video $\mathbf{V}_t \in \mathbb{R}^{F \times C \times H \times W}$ that faithfully follows the target camera trajectory and target intrinsic configuration. Target camera trajectory is defined in a special Euclidean space $(\mathbf{R}, \mathbf{t}) \in SE(3)$, having rotation $\mathbf{R}_t \in \mathbb{R}^{3 \times 3}$ and translation $\mathbf{t}_t \in \mathbb{R}^3$.

In Sec. 4.1, we present our model design tailored for the novel-view video synthesis task. In Sec. 4.2, we de-

scribe our data augmentation strategy that enhances existing synthetic data for novel-view video synthesis across unconstrained camera trajectories and varying intrinsics.

## 4.1. Model Architecture

The overall architecture of our model is illustrated in Fig. 3. We adopt the text-to-video model Wan2.1 [18] as our base architecture. To incorporate camera controllability while preserving Wan2.1's video generation capabilities trained on extensive datasets, we freeze the pretrained weights of Wan2.1 and train only the newly introduced camera encoder and the homography-guided self-attention layers. The new attention layers are initialized using weights from the corresponding pretrained transformer blocks. We employ a camera encoder consisting of a linear layer with 16-dimensional input (comprising a flattened $3 \times 3$ rotation matrix, three translation parameters, focal lengths $f_x$, $f_y$, and principal points $c_x$, $c_y$) and $d$-dimensional output to encode camera poses, where $d$ represents the hidden dimension of the self-attention. For the latent representation $\mathbf{z} \in \mathbb{R}^{b \times f \times h \times w \times d}$, we add the corresponding camera embedding $e_c \in \mathbb{R}^{b \times f \times d}$, broadcasting it along the spatial dimension. To ensure consistent camera conditioning, we share the camera encoder across each DiT block for different camera encodings.

**Warping Module.** Our warping module is illustrated in Fig. 3 (b). The warping module, motivated by Eq. (2), utilizes an infinite homography to warp the latent representation. Since target camera poses are defined relative to the source video's first frame, the warping module warps

the source latent of the initial frame ($\mathbf{z}_s^{\text{init}}$) using the homography derived from the camera rotation and intrinsics, as described in Eq. (1). The warped result is then added to the original $\mathbf{z}_s^{\text{init}}$ through a convolutional layer initialized to zero, functioning as a residual connection. Next, camera embeddings, which encode the target's translation and intrinsics, are incorporated to reflect the second term in Eq. (2). This design simplifies the reprojection estimation under target camera poses by reducing it to parallax relative to $\pi_\infty$ estimation, helping the model achieve higher camera-pose fidelity. The effectiveness of the warping module is validated through ablation studies presented in Table 3.

**Homography-Guided Attention Layer.** The DiT block incorporating the homography-guided attention layer is illustrated in Fig. 3(a). For frame index $i$, we process three types of latents: the source latent $\mathbf{z}_s^i$, target latent $\mathbf{z}_t^i$, and warped latent $\mathbf{z}_w^i$. Each latent is combined with its corresponding camera embedding before being fed into the attention layer. Specifically, to obtain the camera embedding for the target latent, we concatenate the user-specified target intrinsics $\mathbf{K}_t$, rotation $\mathbf{R}_t$, and translation $\mathbf{t}_t$, and feed it into the camera encoder. For the source latent, since the source camera poses are unknown, we form a camera input vector by concatenating the source intrinsics with an identity pose $[\mathbf{I}|\mathbf{0}]$. This vector is then replicated across frames and encoded using the camera encoder. The homography-guided attention layer performs per-frame attention by spatially concatenating these three input latents. The resulting concatenated latents have shape $\mathbf{z}_c \in \mathbb{R}^{bf \times 3hw \times d}$, where frames are processed as individual batch items within the attention mechanism. This structure ensures temporal alignment by allowing the target frames to reference corresponding source frames at the same timestamps. Additionally, the concatenated warped latents help the model better reason about rotation-induced view transformations, which in turn leads to constrained parallax estimation within static regions. After passing through the homography-guided attention layer, the concatenated features are split and reshaped into three tensors of shape $\mathbb{R}^{b \times fhw \times d}$. Finally, we discard the warped latent $\mathbf{z}_i^w$ and feed only the source and target latents into Wan2.1's self-attention layer. During processing through the pretrained Wan2.1 layers, the paired source and target latents are treated as a unified batch. For a more detailed model architecture, please refer to our supplementary material Sec. 7.

## 4.2. Data Preparation

In this section, we describe our data augmentation strategy that enhances existing synthetic data for novel-view video synthesis with unconstrained camera trajectories and varying intrinsics. We utilize the MultiCamVideo Dataset [1] for augmentation and will refer to the augmented version as AugMCV (Augmented MultiCamVideo) Dataset for



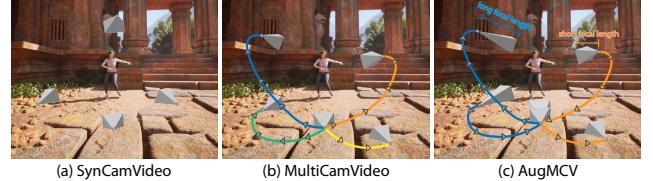(a) SynCamVideo    (b) MultiCamVideo    (c) AugMCV

Figure 4. Visualization of synchronized multi-view synthetic video datasets. Different trajectories are visualized in different colors. (a) **SynCamVideo.** Captured with stationary cameras placed at distinct positions. (b) **MultiCamVideo.** Captured with dynamic cameras following diverse trajectories, all sharing the same initial frame. (c) **AugMCV.** An augmented version of MultiCamVideo with varied starting poses and different focal lengths.

brevity.

### 4.2.1. Trajectory Augmentation

The MultiCamVideo dataset provides 10 trajectories and their corresponding videos for each scene. These trajectories cover a diverse set of camera motions, including pan, tilt, translation, arc, random, and static movements. A notable feature of this data is that all 10 cameras within each scene share identical starting positions, ensuring consistent initialization across viewpoints. Consequently, when two videos from the same scene are randomly sampled as the source and target videos for training, their first frames are always identical. We empirically observe that models trained on this dataset exhibit a bias toward reproducing the source video's first frame, even when conditioned to generate it from a different viewpoint.

Building on these properties, we introduce a data augmentation approach to eliminate such bias from initial frame conditioning. Our key observation is that while all trajectories share the same starting frame, their remaining frames diverge across different videos. Leveraging this property, we randomly sample two distinct videos from the same scene and construct an augmented sequence by reversing the first video and concatenating it with the second one. Since the final frame of the reversed video coincides with the initial frame of the second video, we remove the redundant first frame of the latter, resulting in an augmented video of 161 frames (81 + 80). To align with the pre-training strategy of the Wan2.1 base model, we segment 81 frames from videos within the same scene, ensuring that they share the same starting timestamps for synchronization. Fig. 4 illustrates how the dataset configuration changes under our trajectory augmentation approach. This approach preserves temporal alignment while introducing variability in the initial frame, enabling the model to learn more generalizable trajectory representations and mitigating dataset-specific biases.

### 4.2.2. Intrinsic Augmentation

The MultiCamVideo Dataset incorporates four distinct focal lengths: 18mm, 24mm, 35mm, and 50mm, with the fo-
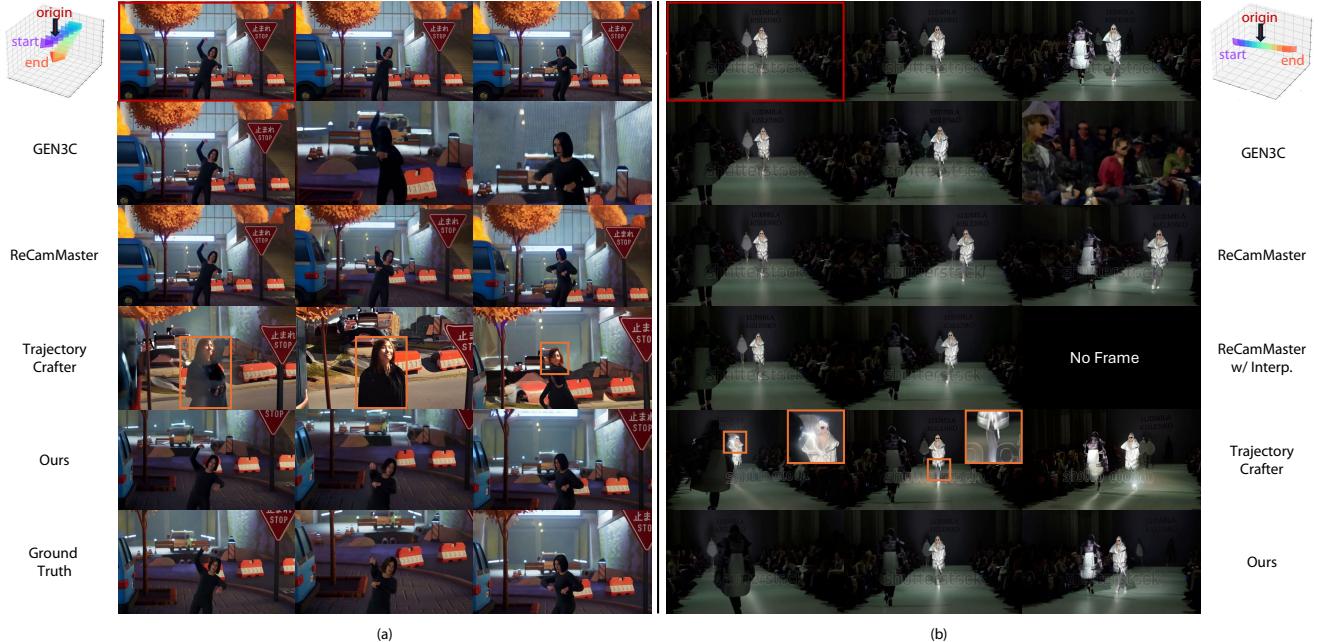
Figure 5. **Qualitative Comparison.** (a) shows results on the test split of the AugMCV dataset, and (b) presents results on the WebVid dataset. In both cases, GEN3C and ReCamMaster fail to perform proper viewpoint transformations, largely preserving the initial frame of the source video. ReCamMaster further fails to reflect pose changes in the initial frame even when trajectory interpolation is applied (ReCamMaster *w/ Interp*). TrajectoryCrafter introduces artifacts due to inaccurate reprojection (highlighted in the orange box). In contrast, our method achieves natural pose transitions while maintaining high visual quality throughout the sequence. Best viewed in zoom.

cal length remaining constant within each scene. Each focal length configuration contains 3,400 scenes, summing up to 13,600 scenes in total. While the dataset inherently encompasses multiple intrinsic parameters, our empirical analysis reveals that models trained on this dataset demonstrate a bias toward generating videos that preserve the focal length of the source video. This limitation arises from the training paradigm wherein both source and target videos are sampled from identical scenes with fixed focal length. Such focal length consistency prevents the model from learning the underlying relationship between focal length variations and their corresponding effects on video generation. To alleviate this bias, we incorporate intrinsic augmentation.

Specifically, given the trajectory-augmented scene with focal length $f_{scene}$, we randomly sample the new focal length $f_{new} \in \{x \in \{18mm, 24mm, 35mm, 50mm\} | x > f_{scene}\}$, and apply intrinsic augmentation process to the input video. Intrinsic augmentation is performed by resizing the video based on the ratio of source and target focal lengths, followed by center cropping. The detailed process is provided in the supplementary material Sec. 8.

#### 4.2.3. Video Pair Selection

To train the video-to-video model for the purpose of novel-view video generation, paired video data is required. In this setup, the source video represents the user input, while the target video serves as the ground truth, guiding the model's output to generate a video that follows the specified trajectory. To create these video pairs, we randomly sample two cameras out of the ten synchronized ones from an identical scene, using one as the source video and the other as the target video. For each video, we apply focal length augmentation with a probability of 0.5. Although we only augment focal lengths in the ascending direction, using augmented videos as either source or target ensures that the selected video pairs encompass both focal length increase and decrease scenarios. This video pair configuration strategy allows for coverage of arbitrary trajectory patterns, as well as the generation of narrow and wide field-of-view videos.

## 5. Experiments

### 5.1. Experiment Settings

**Implementation Details.** We use the pretrained Wan2.1 model as our backbone. Since Wan2.1 requires a text description of the generated video, we employ LLaVA [10] to extract descriptive text from the source video, which we then use as the text input. For training, we use the train split of the AugMCV dataset introduced in Sec. 4.2, which consists of 47,432 scenes, each containing 10 trajectories. For quantitative evaluation, we train on videos of length 81-frame videos at resolution of $480 \times 832$ for 15k iterations using four H100 GPUs (batch size 8), taking approximately one week. For ablations, we train on 41-frame videos at

320 × 544 resolution for 20k iterations using four H100 GPUs (batch size 32), which takes about four days, to ensure efficiency and ease of experimentation. For both resolutions, we use the Adam optimizer with a weight decay of 0.01 and a learning rate of 1e-5.

**Baselines.** We compare our method against state-of-the-art camera-controlled video generation algorithms, including ReCamMaster [1], TrajectoryCrafter [20], and GEN3C [13]. ReCamMaster is a trajectory-conditioned video generation approach trained on synthetic trajectory–video pair datasets. In contrast, TrajectoryCrafter and GEN3C are reprojection-based approaches that estimate depth from the source video and condition the generation process on the reprojection results.

**Evaluation Set.** Our method and baselines are evaluated on two datasets.

*(1) AugMCV Dataset.* Experiments are conducted on the test split of AugMCV dataset, which consists of 168 scenes. Each scene contains one video captured with a static camera and ten videos recorded along dynamic trajectories. Using the static-camera clip as a source, we generate one video for each target trajectory, resulting in a total of 1,680 generated videos. Among the 168 test scenes, 96 scenes have source and target videos with identical focal lengths (referred to as *shared intrinsics*), while the remaining 72 scenes involve different focal lengths between source and target videos (referred to as *different intrinsics*). The generated videos are evaluated against their corresponding ground-truth videos using PSNR, SSIM, and LPIPS [22].

*(2) WebVid Dataset.* A random subset of 100 source videos is selected from the WebVid [3] dataset to evaluate performance in real-world scenarios. For each source video, we generate 20 synthetic videos using different camera trajectories. Specifically, 10 trajectories share the same initial camera pose as the first frame of the source video, while the remaining 10 trajectories start from different initial camera poses. This setup yields a total of 2,000 generated synthetic videos. As the dataset does not provide camera intrinsics for the source videos, we estimate the source intrinsics at inference time using UniDepth [12]. Since no ground-truth videos are available for the target trajectories, we evaluate perceptual fidelity using FID [7] for frame-level quality and FVD [15] for video-level quality. To assess camera pose fidelity, we measure rotation and translation errors [6] between the target trajectories and those estimated from the generated videos using ViPE [8].

## 5.2. Qualitative Results

Fig. 5 (a) and (b) show qualitative results on the AugMCV test split and the WebVid dataset, respectively. In both cases, the first row shows the source video, followed by the generated results from each method. For the AugMCV

Table 1. Quantitative comparison on the AugMCV test set. The best and second-best results are **bold** and underlined, respectively.

| Method | Shared Intrinsics | | | Different Intrinsics | | |
|---|---|---|---|---|---|---|
| | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ |
| GEN3C [13] | 16.891 | 0.479 | 0.548 | 17.449 | 0.525 | 0.467 |
| ReCamMaster [1] | 21.130 | 0.617 | 0.420 | 19.665 | 0.558 | 0.510 |
| TrajectoryCrafter [20] | 21.228 | 0.660 | 0.296 | 19.557 | 0.586 | 0.390 |
| Ours | **22.677** | **0.718** | **0.246** | **22.261** | **0.699** | **0.270** |

Table 2. Quantitative comparison on the WebVid dataset. The best and second-best results are **bold** and underlined, respectively.

| Method | RotErr↓ | TransErr↓ | FID↓ | FVD↓ |
|---|---|---|---|---|
| GEN3C [13] | 9.588 | 3.012 | 43.790 | 331.768 |
| ReCamMaster [1] | 8.375 | 1.027 | 39.930 | 302.025 |
| ReCamMaster w/ interp. | 9.016 | 1.590 | 35.720 | 303.012 |
| TrajectoryCrafter [20] | 5.007 | 0.735 | 30.877 | 289.879 |
| Ours | **3.162** | **0.438** | **29.702** | **286.952** |

dataset, the last row displays the ground-truth target video, whereas for the WebVid dataset, only the generated results are available since ground-truth target views are not provided. Frames are uniformly sampled along the camera trajectory for visualization.

Across both datasets, ReCamMaster consistently preserves the initial frame of the source video, which we attribute to a bias in its training data, where all paired source-target trajectories originate from the same initial frame of the source video. To mitigate this bias, we applied a frame interpolation strategy and visualized the results in the fourth row of Fig. 5(b). Specifically, eight identical frames were prepended to the source video, and their corresponding camera poses were linearly interpolated between the identity pose and the target pose of the first frame. The first eight frames were then removed from the generated video before evaluation to exclude the interpolated portion. Although this strategy enables ReCamMaster to synthesize different viewpoints at the initial frame, the results of *ReCamMaster w/ interp.* in Fig. 5(b) still exhibit alignment errors, indicating that interpolation alone cannot fully address its trajectory bias.

TrajectoryCrafter, a reprojection-based method, successfully performs viewpoint transformation but fails to maintain the source appearance throughout the sequence due to reprojection errors introduced by inaccurate depth estimation. GEN3C also relies on reprojection, but fails to transform the viewpoint of the initial frame. We conjecture that this result stems from a bias in its image-to-video backbone toward preserving the initial source frame. This observation is consistent with that reported in GCD [16].

In contrast, our method, empowered by the infinite-homography-based warping module and trajectory-intrinsic augmentation, achieves consistent viewpoint alignment and maintains high visual coherence throughout the trajectory.

Table 3. Ablation study conducted on the AugMCV testset. The best and second-best results are **bold** and underlined, respectively.

| Components | | | Shared Intrinsic | | | Different Intrinsic | | |
|---|---|---|---|---|---|---|---|---|
| Aug.Traj. | Aug.Intr. | Warp | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ |
| □ | □ | □ | 19.228 | 0.562 | 0.427 | 18.480 | 0.507 | 0.525 |
| ✓ | □ | □ | 20.820 | 0.615 | 0.353 | 18.865 | 0.523 | 0.499 |
| ✓ | ✓ | □ | <u>22.807</u> | <u>0.680</u> | <u>0.250</u> | <u>21.866</u> | <u>0.649</u> | <u>0.293</u> |
| ✓ | ✓ | ✓ | **24.412** | **0.733** | **0.198** | **24.311** | **0.720** | **0.203** |

Table 4. Ablation study of our augmentation strategies on WebVid dataset. Best scores per metric are in **bold**.

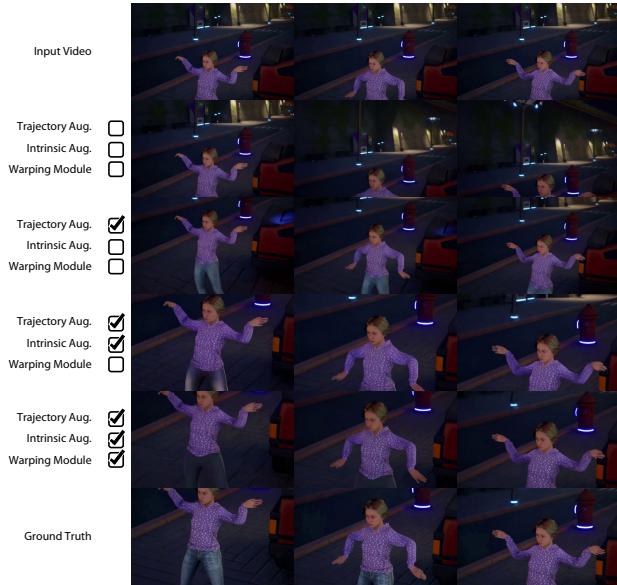| Method | Rot.↓ | Trans.↓ | FID↓ | FVD↓ |
|---|---|---|---|---|
| *w/* MCV+SCV | 4.158 | 1.441 | 47.949 | 235.958 |
| *w/* AugMCV | **3.368** | **0.839** | **40.384** | **235.563** |



Figure 6. Qualitative ablation study of proposed components. From top to bottom, each row incrementally adds proposed components, showing cumulative improvements toward the target video. Best viewed in zoom.

## 5.3. Quantitative Results

Table 1 presents quantitative evaluation results on AugMCV dataset for two scenarios: (1) source and target videos with identical camera intrinsics, and (2) source and target videos with different camera intrinsics. Across all scenarios, our method consistently outperforms the baselines in all metrics, demonstrating that the videos generated by our method are clearly closer to the ground truth.

Table 2 presents the quantitative evaluation results on the WebVid dataset. Our method consistently outperforms the baseline approaches in terms of both camera pose accuracy and visual fidelity. In particular, camera pose accuracy demonstrates a significant improvement over the baselines.

## 5.4. Ablation Study

**Analysis of Proposed Components.** We perform an ablation study on AugMCV test set by progressively adding each of the proposed components. As shown in Table 3 and Fig. 6, the baseline without warping module and data augmentation fails to capture both the target trajectory and intrinsics, yielding the lowest scores in all settings. Adding trajectory augmentation improves all scores and results in

a movement consistent with the target trajectory in Fig. 6; however, it still fails to reflect the target intrinsics, offering comparable performance to the baseline under *different intrinsics* setting. Combining intrinsic augmentation with trajectory augmentation allows the model to better understand the different intrinsic settings, though performance remains inaccurate when compared to the ground truth. Finally, introducing the warping module, which explicitly warps the source latent using the target pose and intrinsics, proves crucial for accurately aligning the results with the ground truth, leading to substantial qualitative and quantitative improvements.

**Effectiveness of our augmentation strategy.** We evaluate the proposed trajectory-intrinsic augmentation scheme by comparing our model trained on the AugMCV dataset with a model trained on the mixture of the MultiCamVideo (MCV) dataset and the SynCamVideo (SCV) dataset. The SynCamVideo dataset is a synchronized multi-camera video dataset rendered using Unreal Engine 5, where the cameras are stationary and their corresponding poses are provided. Unlike the MultiCamVideo dataset, the SynCamVideo dataset has different initial frame poses for the source and target videos. Jointly training on both datasets is expected to mitigate the bias inherent in each dataset. However, as shown in Table 4, our model trained with AugMCV outperforms the joint training approach, highlighting the superior effectiveness of our trajectory-intrinsic augmentation.

## 6. Conclusion

In this work, we introduced InfCam, a depth-free framework that achieves precise camera control in video generation through infinite homography warping. By encoding camera rotations in the latent space and learning residual parallax, our method removes the dependence on external depth estimation modules and erroneous reprojection results, resulting in a significant enhancement in camera-pose accuracy. Combined with our trajectory-intrinsic augmentation strategy, InfCam demonstrates superior performance across diverse camera motions and focal lengths. InfCam consistently outperforms depth reprojection-based and trajectory-conditioned approaches across all metrics. Future work could involve extending this framework to longer video sequences, allowing for comprehensive camera-controlled generation across extended temporal horizons.

# References

[1] Jianhong Bai, Menghan Xia, Xiao Fu, Xintao Wang, Lianrui Mu, Jinwen Cao, Zuozhu Liu, Haoji Hu, Xiang Bai, Pengfei Wan, et al. Recammaster: Camera-controlled generative rendering from a single video. In *ICCV*, 2025. 2, 3, 5, 7, 1

[2] 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. *The Thirteenth International Conference on Learning Representations*, 2025. 2

[3] Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. Frozen in time: A joint video and image encoder for end-to-end retrieval. In *IEEE International Conference on Computer Vision*, 2021. 7

[4] 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 *Forty-first international conference on machine learning*, 2024. 3

[5] Richard Hartley and Andrew Zisserman. Scene planes and homographies. In *Multiple view geometry in computer vision*, pages 325–342. Cambridge university press, 2003. 3

[6] Hao He, Yinghao Xu, Yuwei Guo, Gordon Wetzstein, Bo Dai, Hongsheng Li, and Ceyuan Yang. Cameractrl: Enabling camera control for video diffusion models. In *The Thirteenth International Conference on Learning Representations*, 2025. 7, 3

[7] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017. 7

[8] Jiahui Huang, Qunjie Zhou, Hesam Rabeti, Aleksandr Korovko, Huan Ling, Xuanchi Ren, Tianchang Shen, Jun Gao, Dmitry Slepichev, Chen-Hsuan Lin, Jiawei Ren, Kevin Xie, Joydeep Biswas, Laura Leal-Taixe, and Sanja Fidler. Vipe: Video pose engine for 3d geometric perception. In *NVIDIA Research Whitepapers*, 2025. 7, 3

[9] Hanwen Liang, Junli Cao, Vidit Goel, Guocheng Qian, Sergei Korolev, Demetri Terzopoulos, Konstantinos N Plataniotis, Sergey Tulyakov, and Jian Ren. Wonderland: Navigating 3d scenes from a single image. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 798–810, 2025. 2

[10] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, 2024. 6

[11] Norman Müller, Katja Schwarz, Barbara Rössle, Lorenzo Porzi, Samuel Rota Bulo, Matthias Nießner, and Peter Kontschieder. Multidiff: Consistent novel view synthesis from a single image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10258–10268, 2024. 2

[12] Luigi Piccinelli, Yung-Hsu Yang, Christos Sakaridis, Mattia Segu, Siyuan Li, Luc Van Gool, and Fisher Yu. Unidepth: Universal monocular metric depth estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10106–10116, 2024. 7

[13] 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 *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 6121–6132, 2025. 2, 7, 3

[14] Wenqiang Sun, Shuo Chen, Fangfu Liu, Zilong Chen, Yueqi Duan, Jun Zhu, Jun Zhang, and Yikai Wang. Dimensionx: Create any 3d and 4d scenes from a single image with decoupled video diffusion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13695–13706, 2025. 2

[15] Thomas Unterthiner, Sjoerd Van Steenkiste, Karol Kurach, Raphaël Marinier, Marcin Michalski, and Sylvain Gelly. Fvd: A new metric for video generation. 2019. 7

[16] 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. In *European Conference on Computer Vision*, pages 313–331. Springer, 2024. 2, 7

[17] Vikram Voleti, Chun-Han Yao, Mark Boss, Adam Letts, David Pankratz, Dmitry 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*, pages 439–457. Springer, 2024. 2

[18] 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, 2

[19] Dejia Xu, Weili Nie, Chao Liu, Sifei Liu, Jan Kautz, Zhangyang Wang, and Arash Vahdat. Camco: Camera-controllable 3d-consistent image-to-video generation. *arXiv preprint arXiv:2406.02509*, 2024. 2

[20] Mark YU, Wenbo Hu, Jinbo Xing, and Ying Shan. Trajectorycrafter: Redirecting camera trajectory for monocular videos via diffusion models. In *ICCV*, 2025. 2, 7, 3

[21] 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 *Proceedings of the Computer*

*Vision and Pattern Recognition Conference*, pages 2050–2062, 2025. 2

[22] 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

# Infinite-Homography as Robust Conditioning for Camera-Controlled Video Generation
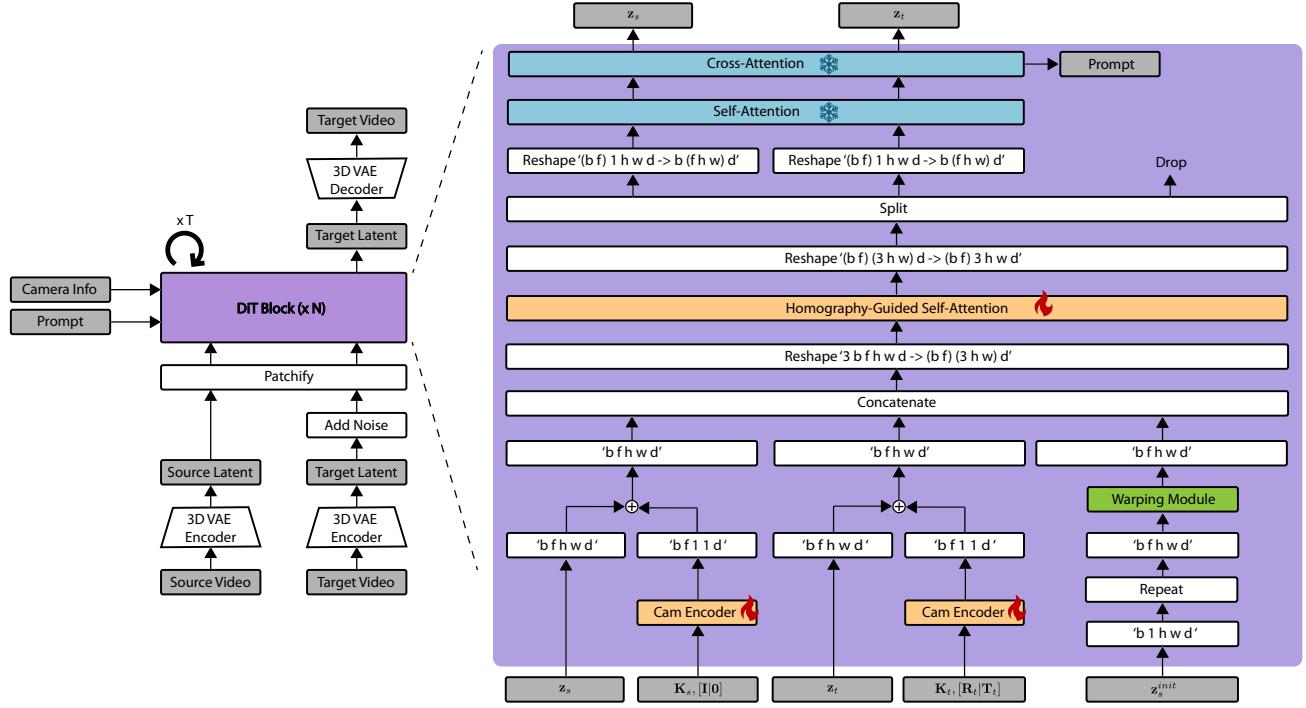
## Supplementary Material



Figure 7. This figure illustrates our model overview along with the dimensions of latent features. Here, $b$ denotes the batch size, $f$ indicates the latent frame count, $w$ and $h$ represent latent width and height, and $d$ represents the feature dimension. The dimensional notation follows the einops convention.

## 7. Model Architecture Details

Fig. 7 illustrates the overall architecture of the proposed InfCam with its dimensional specifications. As described in Sec. 4.1, the homography-guided self-attention layer performs per-frame attention by spatially concatenating three input latents: the source, target, and warped latents. The resulting concatenated latent $\mathbf{z}_c \in \mathbb{R}^{bf \times 3hw \times d}$ treats frames as individual batch items within the attention mechanism, where spatial features from all three latents are concatenated along the spatial dimension.

After passing through the homography-guided self-attention layer, $\mathbf{z}_c$ is split and reshaped back into three separate tensors, each with shape $\mathbb{R}^{b \times fhw \times d}$. Finally, we discard the warped latent and feed only the source and target latents into Wan2.1's self-attention layer. During processing through the pretrained Wan2.1 layers, the paired source and target latents are treated as a unified batch.

## 8. Data Preparation Details

**Trajectory Augmentation.** Fig. 8 (a) shows an example scene from the MultiCamVideo dataset [1]. The MultiCamVideo dataset provides 10 trajectories and their corresponding videos for each scene. These trajectories cover a diverse set of camera motions, including pan, tilt, translation, arc, random, and static movements. A notable property of this dataset is that all 10 cameras within each scene share an identical starting position (highlighted in yellow), which ensures that the first frame of the target video is always identical to the first frame of the source video. We hypothesize that these dataset constraints arise from the fundamental requirement that camera viewpoints must maintain sufficient frustum overlap for the purpose of novel-view video generation, ensuring the task does not collapse into pure video generation. To reduce this bias, we perform the data augmentation introduced in Sec. 4.2. Fig. 8 (b) illustrates an augmented trajectory example and its correspond-
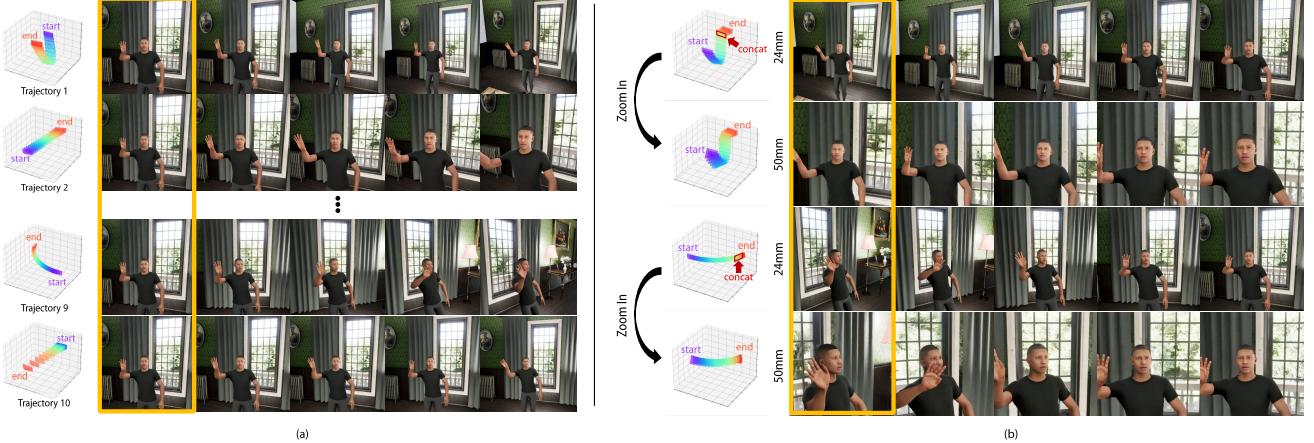
Figure 8. **(a)** Example scene from the MultiCamVideo dataset. Each row shows a camera trajectory and its corresponding video, with all videos sharing the same initial frame (highlighted in yellow). **(b)** Augmented dataset examples. Rows 1&3 show trajectory augmentation, which preserves temporal alignment among trajectories while introducing variation in the initial frame selection (highlighted in yellow). Rows 2&4 show focal-length augmentation, enabling the model to learn the relationship between focal-length changes and their visual effects in video generation.

ing focal-length (intrinsics) augmentation. The trajectories illustrated in the first and third rows of Fig. 8 (b) are constructed from the trajectory pairs (Trajectory 1, Trajectory 2) and (Trajectory 9, Trajectory 10) in Fig. 8 (a) by reversing the first trajectory and concatenating it with the second one. From the resulting sequence, a segment of length 81 frames is randomly cropped to match the supported sequence length of the Wan2.1 base model.

**Intrinsic Augmentation.** Intrinsic augmentation is applied to the augmented trajectories. For illustration, the first and third rows in Fig. 8 (b) depict the trajectories prior to intrinsic augmentation, whereas the second and fourth rows display their augmented counterparts. Algorithm 1 illustrates how we augment the focal length for a given scene. This process involves resizing according to the ratio between the current and new focal lengths, followed by center cropping to maintain the original image resolution. Since we only augment focal lengths in the ascending direction as stated in Sec. 4.2.2, the resized and cropped images remain within the originally observed field of view and do not introduce any unseen regions.

**Video Pair Selection.** During training, given a trajectory pair and its corresponding intrinsically augmented pair as in Fig. 8 (b), we randomly select two trajectories among these four trajectories, using one as the source and the other as the target. Although we only augment focal lengths in the ascending direction, using augmented videos as either source or target ensures that the selected video pairs encompass both focal-length increase and decrease scenarios.

**Overall.** By applying augmentations, we expand the dataset from 13,600 to a total of 47,600 videos. Among these, 47,432 videos are used for training, and 168 videos

are reserved for testing, without scene overlap between the sets.

---

**Algorithm 1** Intrinsic Augmentation

---

1: **Input:** Input/Output scene paths $P_{scene}^{in}$, $P_{scene}^{out}$, focal lengths $f_{old}, f_{new}$
2: **Output:** Augmented video dataset with updated focal length
3: **procedure** PROCESSVIDEOS($P_{scene}^{in}$, $P_{scene}^{out}$, $f_{now}$, $f_{new}$)
4:     **for** $i = 1$ to 10 **do**                ▷ Process 10 cameras
5:         $video \leftarrow$ Load($P_{scene}^{in}$/cam_$i$.mp4)
6:         $(W_{orig}, H_{orig}) \leftarrow$ GetDimensions($video$)
7:         $(W_{new}, H_{new}) \leftarrow (\frac{f_{new}}{f_{old}} \cdot W_{orig}, \frac{f_{new}}{f_{old}} \cdot H_{orig})$
8:         **while** frame exists in $video$ **do**
9:             $frame \leftarrow$ ReadFrame($video$)
10:             $frame \leftarrow$ Resize($frame$, $(W_{new}, H_{new})$)
11:             $frame \leftarrow$ CenterCrop($frame$, $(W_{orig}, H_{orig})$)
12:             WriteFrame($P_{scene}^{out}$/cam_$i$.mp4, $frame$)
13:         **end while**
14:     **end for**
15: **end procedure**

---

## 9. Additional Experimental Details

### 9.1. Loss Function

We employ the same training objective as Wan2.1 [18].

**Wan2.1** [18] Wan2.1 is an open-source Text-to-Video (T2V) diffusion model based on a transformer architecture. During training, for a given video $\mathbf{V} \in \mathbb{R}^{B \times (1+F) \times H \times W \times 3}$, the Wan-VAE compresses its spatio-temporal dimensions from $(1 + F, H, W)$ to $(1 + F/4, H/8, W/8)$. Subsequent patchification further reduces the spatial resolution, yielding $\mathbf{z} \in \mathbb{R}^{B \times (f \times h \times w) \times d}$

, where $f = 1 + F/4$, $h = H/16$, $w = W/16$. Given a video latent $\mathbf{z}_1$, a random noise $\mathbf{z}_0 \sim \mathcal{N}(0, I)$, and a sampled timestep $t \in [0, 1]$, an intermediate latent $\mathbf{z}_t$ is obtained as the training input. Following Rectified Flows [4], $\mathbf{z}_t$ is defined as a linear interpolation between $\mathbf{z}_0$ and $\mathbf{z}_1$, i.e., $\mathbf{z}_t = t\mathbf{z}_1 + (1 - t)\mathbf{z}_0$. The ground truth velocity $\mathbf{v}_t$ is $\mathbf{v}_t = \frac{d\mathbf{z}_t}{dt} = \mathbf{z}_1 - \mathbf{z}_0$. The model is trained to predict the velocity, thus, the loss function can be formulated as the mean squared error (MSE) between the model output and $\mathbf{v}_t$,

$$\mathcal{L} = \mathbb{E}_{\mathbf{z}_0, \mathbf{z}_1, c_{txt}, t} \left[ \| u(\mathbf{z}_t, c_{txt}, t; \theta) - \mathbf{v}_t \|^2 \right], \quad (3)$$

where $c_{txt}$ is the text embedding sequence, $\theta$ represents the model weights, and $u(\mathbf{z}_t, c_{txt}, t; \theta)$ denotes the output velocity predicted by the model.

### 9.2. Baseline Setup

**ReCamMaster** [1]. Following the official code and checkpoints, we generate videos at a resolution of $480 \times 832$ and a length of 81 frames. In addition, as described in Sec. 5.2, we introduce *ReCamMaster w/ interp.*, which mitigates ReCamMaster's bias to preserve the first frame of the source video by prepending eight auxiliary frames that smoothly interpolate the camera pose of the first frame toward the target pose.

**TrajectoryCrafter** [20]. Following the official code and checkpoints, we generate videos at a resolution of $384 \times 672$ and a length of 49 frames. For comparison under the 81-frame setting, we first generate a 49-frame segment, and then generate an additional 33-frame segment. When generating the second segment, we feed the last frame of the first segment output as the first frame of the reprojection condition for the second segment to enforce temporal consistency at the segment boundary. After inference, we discard the first frame of the second segment and concatenate the remaining 32 frames to the initial 49-frame segment, obtaining an 81-frame sequence with a one-frame overlap removed. This inference-level extension alleviates the architectural limitation while preserving temporal continuity, and all other experimental settings are kept identical to those in the original paper.

**GEN3C** [13]. Following the official code and checkpoints, we generate videos at a resolution of $704 \times 1280$ and a length of 121 frames. For the 81-frame comparison, we extend the input sequence by duplicating its final frame 40 times and appending these frames to the end of the sequence. GEN3C is conditioned on this extended input, and the first 81 frames of the generated output are extracted and used for evaluation.

### 9.3. Camera Accuracy Metric

For RotErr (degree) and TransErr (meter) computation, ViPE [8] is used for camera trajectory extraction of the generated video, and the extracted trajectory is compared with the ground truth trajectory. To ensure that camera trajectory extraction with ViPE [8] references the initial camera pose of the source video, the first frame of each source video is prepended to every generated video. ViPE then estimates relative poses with respect to this concatenated first frame. The pose estimated for this prepended frame is discarded before evaluation. Following [6], we compute the rotation error and translation error as

$$RotErr = \arccos\left( \left( \mathrm{tr}\left( \mathbf{R}_{\mathrm{pred}} \mathbf{R}_{\mathrm{gt}}^\top \right) - 1 \right) / 2 \right), \quad (4)$$

$$TransErr = \| \mathbf{t}_{\mathrm{pred}} - \mathbf{t}_{\mathrm{gt}} \|_2, \quad (5)$$

where $\mathbf{R}_{\mathrm{pred}}, \mathbf{R}_{\mathrm{gt}} \in \mathrm{SO}(3)$ are the predicted and ground-truth rotation matrices, respectively, and $\mathbf{t}_{\mathrm{pred}}, \mathbf{t}_{\mathrm{gt}} \in \mathbb{R}^3$ are their corresponding translation vectors.

## 10. Additional Results

### 10.1. Additional Quantitative Results

Table 5 and Table 6 present the quantitative results for the First-Frame Synchronized (FF-Sync) and First-Frame Asynchronous (FF-Async) settings. FF-Sync refers to the experiment setting where the source and target videos share the same initial frame, whereas FF-Async uses source and target videos with different initial frames, yielding a significantly more challenging evaluation setting. Our method demonstrates clear superiority over competing approaches in the FF-Async setting and achieves competitive or superior performance in the FF-Sync setting. In particular, our method attains lower rotation and translation errors in the more challenging FF-Async scenario. These results validate that the proposed warping and augmentation strategies

Table 5. Quantitative comparison of FF-Sync case on Web-Vid dataset. The best results are in **bold**; the second-best are underlined.

| Method | FF-Sync | | | |
| | RotErr↓ | TransErr↓ | FID↓ | FVD↓ |
| --- | --- | --- | --- | --- |
| GEN3C [13] | 10.012 | 2.329 | 53.103 | 350.499 |
| ReCamMaster [1] | 9.673 | 1.466 | 40.612 | 308.697 |
| TrajectoryCrafter [20] | <u>5.595</u> | **0.502** | **32.220** | <u>287.805</u> |
| Ours | **3.605** | <u>0.510</u> | <u>32.906</u> | **282.703** |

Table 6. Quantitative comparison of FF-Async case on Web-Vid dataset. The best results are in **bold**; the second-best are underlined.

| Method | FF-Async | | | |
| | RotErr↓ | TransErr↓ | FID↓ | FVD↓ |
| --- | --- | --- | --- | --- |
| GEN3C [13] | 9.165 | 3.694 | 34.476 | 313.036 |
| ReCamMaster [1] | 7.076 | <u>0.589</u> | 39.248 | <u>295.353</u> |
| ReCamMaster w/ interp. | 8.360 | 1.714 | <u>30.828</u> | 297.328 |
| TrajectoryCrafter [20] | <u>4.418</u> | 0.969 | 34.476 | 313.036 |
| Ours | **2.718** | **0.365** | **26.497** | 291.202 |

effectively generate videos that remain well-aligned with the target camera trajectories, regardless of the difficulty of the given trajectory configuration. To compute the results in Table 2, we average the scores from Table 5 and Table 6. When performing this averaging, the value corresponding to *ReCamMaster w/ interp.* in Table 6 is paired with the value of the original ReCamMaster (without interpolation) from Table 5, since interpolation is unnecessary in the synchronized case.

## 10.2. Additional Qualitative Results

Fig. 9 and Fig. 10 provide additional qualitative results on both the AugMCV and WebVid datasets. Fig. 11 further illustrates the generalized performance of our method across diverse camera trajectories. All corresponding video results are included in the supplementary material as separate video files.
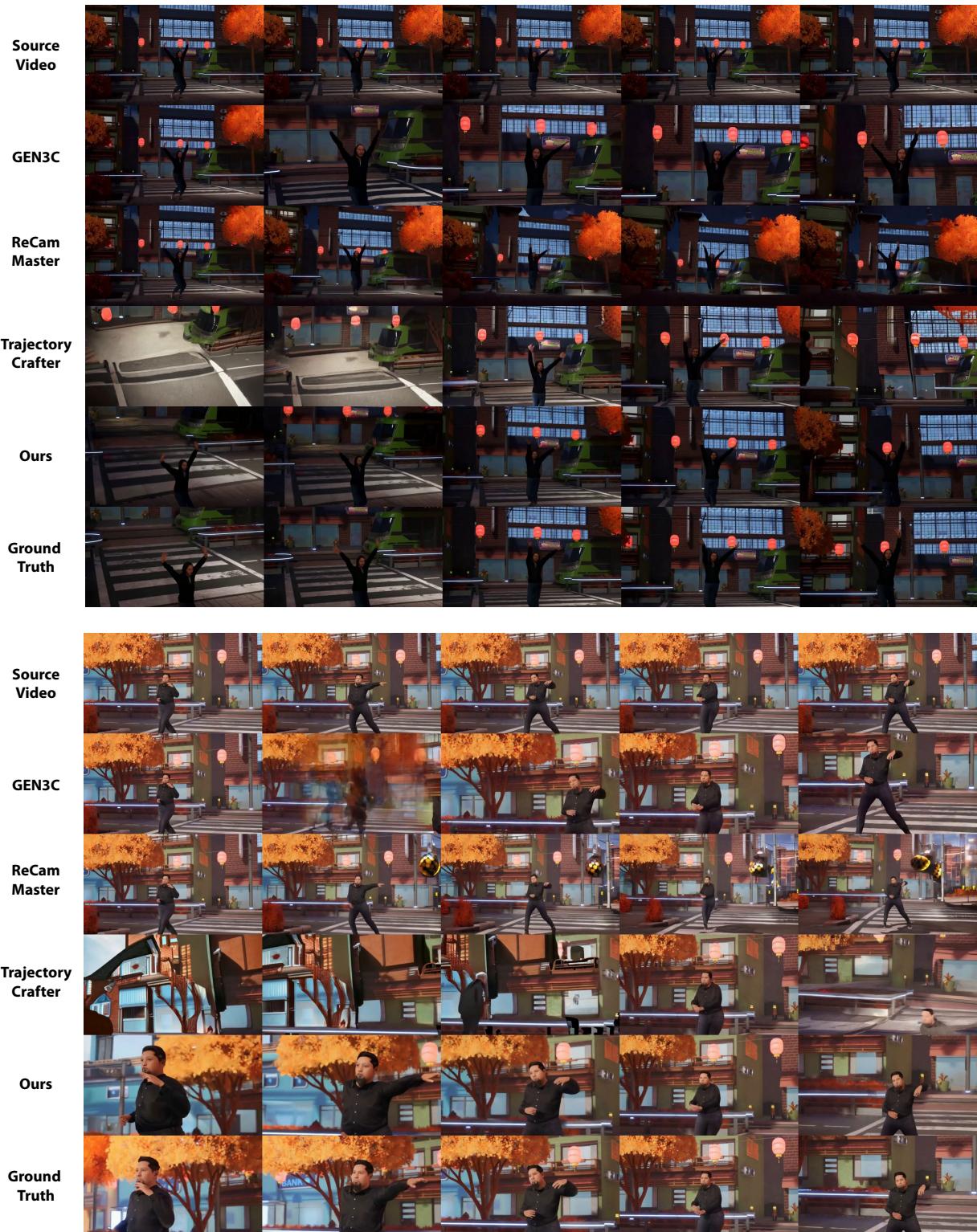
Figure 9. Additional qualitative results on AugMCV dataset. Best viewed in zoom.

Figure 10. Additional qualitative results on WebVid dataset. The first row illustrates the target trajectory and source video, while the remaining rows present the outputs of the baseline methods and our method. Note that both GEN3C and ReCamMaster fail to modify the camera pose of the first frame (highlighted in yellow). For ease of comparison, the key regions for comparison are highlighted in orange. Best viewed in zoom.

Figure 11. Additional qualitative results of our method under eight different camera trajectories. Best viewed in zoom.