3D SCENE PROMPTING FOR SCENE-CONSISTENT CAMERA-CONTROLLABLE VIDEO GENERATION

Anonymous authors

000

001

002003004

010 011

012

013

014

016

017

018

019

021

024

025

026

027

028

029

031

033

035

037

040

041

042

043

044

046 047

048

051

052

Paper under double-blind review

ABSTRACT

We present **3DScenePrompt**, a framework for camera-controllable video generation that maintains scene consistency when extending arbitrary-length input videos along user-specified trajectories. Unlike existing video generative methods limited to conditioning on a single image or just a few frames, we introduce a dual spatio-temporal conditioning strategy that fundamentally rethinks how video models should reference prior content. Our approach conditions on both temporally adjacent frames for motion continuity and spatially adjacent content for scene consistency. However, when generating beyond temporal boundaries, directly using spatially adjacent frames would incorrectly preserve dynamic elements from the past. We address this through introducing a 3D scene memory that represents exclusively the static geometry extracted from the entire input video. To construct this memory, we leverage dynamic SLAM with our newly introduced dynamic masking strategy that explicitly separates static scene geometry from moving elements. The static scene representation can then be projected to any target viewpoint, providing geometrically-consistent warped views that serve as strong spatial prompts while allowing dynamic regions to evolve naturally from temporal context. This enables our model to maintain long-range spatial coherence and precise camera control without sacrificing computational efficiency or motion realism. Extensive experiments demonstrate that our framework significantly outperforms existing methods in scene consistency, camera controllability, and generation quality.

1 Introduction

Camera-controllable video generation (He et al., 2024; Wang et al., 2024b; Jin et al., 2025) aims to synthesize videos following user-specified camera trajectories while maintaining visual coherence and temporal consistency. Recent advances have progressed from generating entirely new videos with controllable viewpoints (Bahmani et al., 2025) to enabling users to extend a single image or short video clips along desired camera paths (He et al., 2024; Agarwal et al., 2025). Yet these methods share a fundamental limitation: they can only process extremely short conditioning sequences, typically just a few frames, which constrains their ability to understand longer videos and hence fails to preserve the rich scene context present in those longer videos. What if we could provide a model with arbitrary-length video sequences and generate continuations that not only follow precise camera controls but also maintain scene consistency with the entire input? Such technology, which we refer to as scene-consistent camera-controllable video generation, has immediate applications in film production (Zhang et al., 2025), virtual reality (He et al., 2025), and synthetic data generation (Knapp & Bohacek, 2025).

Scene-consistent camera-controllable video generation poses three intertwined challenges that must be solved jointly. First, static and dynamic elements must be handled differently: while static scene elements should remain consistent throughout generation, dynamic elements such as moving objects and people should evolve naturally from their most recent states rather than rigidly preserving motions from the distant past. Second, camera control demands understanding the underlying 3D geometry of the scene: the generated content must respect physical constraints, properly handle occlusions, and seamlessly compose dynamic elements onto static geometry, while extrapolating plausible content for previously unobserved regions. Third, these capabilities must be achieved

within practical computational constraints, as naive approaches that process all input frames quickly become intractable when the input video sequence is long.

How can we tackle this challenging task by leveraging existing video generative models? Our key insight comes from fundamentally rethinking how video models should reference prior content. Current image-to-video (Yang et al., 2024) and video-to-future-video models ¹ (Agarwal et al., 2025) achieve realistic generation by conditioning on temporally adjacent frames to maintain short-term consistency and motion continuity. However, adjacency in video is not purely temporal—it can also be spatial. When generating scene-consistent videos, the frames we synthesize may be spatially adjacent to frames from much earlier in the input sequence, particularly when the camera revisits similar viewpoints or explores nearby regions. This dual nature of adjacency suggests a new conditioning paradigm that leverages both temporal and spatial relationships.

Based on these motivations, we propose **3DScenePrompt**, a novel video generation framework designed for scene-consistent camera-controllable video synthesis. It takes an arbitrary-length video as context and generates the future video that is consistent with the scene geometry of the context video. The key innovation lies in our dual spatio-temporal conditioning strategy: the model conditions on both *temporally* adjacent frames (for motion continuity) and *spatially* adjacent frames (for scene consistency). However, an important consideration for spatial conditioning for our task is that it must provide only the persistent *static* scene structure while excluding *dynamic* content, as directly conditioning on spatially adjacent frames from the past would incorrectly preserve dynamic elements. To enable this without temporal contradictions, we construct a **3D scene memory** that represents exclusively the *static* geometry extracted from the entire input video.

To construct this 3D scene memory from *dynamic* videos, we leverage recent advances in dynamic SLAM frameworks (Zhang et al., 2022; 2024; Li et al., 2024) to estimate camera poses and 3D structure from the input video. To extract only the *static* regions from the estimated 3D structure, we introduce a dynamic masking strategy that explicitly separates static elements and moving objects. The static-only 3D representation can then be projected to target viewpoints, yielding geometrically-consistent warped views that serve as *spatial prompts* while allowing dynamic elements to evolve naturally from temporal context alone. Surprisingly, the integration of 3D scene memory provides an additional benefit: the geometrically-consistent warped views provide rich visual references that significantly reduce uncertainty in viewpoint manipulation, enabling precise camera control without any other explicit camera conditioning.

In summary, **3DScenePrompt** enables both accurate camera control and long-range spatial consistency by treating the static scene representation as a persistent spatial prompt that guides generation across arbitrary timescales. Extensive experiments demonstrate that our framework significantly outperforms existing methods in maintaining scene consistency, achieving precise camera control, and generating high-quality videos from arbitrary-length inputs.

2 RELATED WORK

Camera-controllable video generation. Building upon the recent success of video diffusion models (Blattmann et al., 2023; Guo et al., 2023; Yang et al., 2024; Runway; Brooks et al., 2024), recent works (He et al., 2024; Wang et al., 2024b; Bahmani et al., 2024) have achieved camera-controllable video generation by introducing additional adapters into U-Net-based video diffusion models that accept camera trajectories. For instance, CameraCtrl and VD3D (Bahmani et al., 2024; He et al., 2024) incorporate spatiotemporal camera embeddings, such as Plücker coordinates, via ControlNet-like mechanisms (Zhang et al., 2023). While these methods enable precise trajectory following, they only condition on single starting images, lacking mechanisms to maintain consistency with extended video context. In contrast, our approach enables leveraging entire video sequences as spatial prompts through 3D memory construction, enabling scene-consistent generation that preserves the rich scene context within arbitrary-length inputs.

Geometry-grounded video generation. Recent works (Ren et al., 2025; Yu et al., 2025; Seo et al., 2025) have integrated off-the-shelf geometry estimators into video generation pipelines to improve

¹Throughout our paper, video-to-future-video models refer to models that are capable of generating the subsequent frames of the given input video (e.g., cosmos-predict2 (Agarwal et al., 2025).

geometric accuracy. Gen3C (Ren et al., 2025), for instance, similarly adopts dynamic SLAM to lift videos to 3D representations. However, these methods exclusively address dynamic novel view synthesis—generating new viewpoints within the same temporal window as the input. This constrained setting allows them to simply warp entire scenes without distinguishing static and dynamic elements. Our work fundamentally differs by generating content beyond temporal boundaries, requiring selective masking of dynamic regions during 3D construction—a critical challenge that emerges only when static geometry must persist while dynamics evolve naturally into the future.

Long-horizon scene-consistent generation. Various approaches attempt scene-consistent long video generation through different strategies. ReCamMaster (Bai et al., 2025) and TrajectoryCrafter (Yu et al., 2025) interpolate frames or construct 3D representations but remain confined to the input's spatiotemporal coverage, essentially performing dynamic novel view synthesis. Star-Gen (Zhai et al., 2025) scales to long trajectories but assumes static worlds, eliminating temporal dynamics entirely. DFoT (Song et al., 2025) most closely relates to our work, proposing guidance methods that condition on previous frames for scene consistency. However, DFoT also faces fundamental memory constraints when processing extended sequences, limiting its ability to maintain long-range spatial coherence. Our dual spatio-temporal strategy with SLAM-based spatial memory overcomes these limitations by selectively retrieving only the most relevant frames, both temporally and spatially, enabling computationally efficient processing of arbitrary-length videos while maintaining both motion continuity and scene consistency.

3 METHODOLOGY

3.1 PROBLEM FORMULATION AND MOTIVATION

We address the task of scene-consistent camera-controllable video generation: given a dynamic video $\mathbf{V}_{\text{in}} \in \mathbb{R}^{L \times H \times W \times 3}$ of arbitrary length L as context with height H and width W, our goal is to generate T subsequent frames $\mathbf{V}_{\text{out}} \in \mathbb{R}^{T \times H \times W \times 3}$ that follow a desired camera trajectory $\mathbf{C} = \{C_t\}_{t=1}^T$ while maintaining consistency with the scene captured in the context input:

$$\mathbf{V}_{\text{out}} = \mathcal{F}(\mathbf{V}_{\text{in}}, \mathcal{T}, \mathbf{C}), \tag{1}$$

where $C_t \in \mathbb{SE}(3)$ represents camera extrinsic matrices and \mathcal{T} is a text prompt when a video generator $\mathcal{F}(\cdot)$ is based on pretrained text-to-video priors (Yang et al., 2024; Bahmani et al., 2025).

Comparison to existing solutions. This task fundamentally differs from existing video generation paradigms. Existing camera-controllable generation methods (He et al., 2024; Wang et al., 2024b; Bahmani et al., 2024) synthesize videos following user-specified trajectories but only condition on a single image I_{ref} or plain text \mathcal{T} (Fig. 1-(a)):

$$\mathbf{V}_{out} = \mathcal{F}(\mathbf{I}_{ref}, \mathcal{T}, \mathbf{C}), \text{ or } \mathbf{V}_{out} = \mathcal{F}(\mathcal{T}, \mathbf{C}),$$
 (2)

which is insufficient for our task, where the entire underlying 3D scene of the context video should be considered. In contrast, video-to-future-video generation methods such as Cosmos-predict-2 (Agarwal et al., 2025) $\mathcal{G}(\cdot)$ employ temporal sliding windows to generate future frames (Fig. 1-(b)):

$$\mathbf{V}_{\text{out}} = \mathcal{G}(\mathbf{V}_{\text{in}}[L - w : L], \mathcal{T})$$
(3)

where $V_{\rm in}[L-w:L]$ for $w\ll L$ represents a small overlap window, typically consisting of the last few frames of $V_{\rm in}$. Although this design encourages temporal smoothness by providing the last few frames when generating the future video, it often fails to preserve long-term spatial consistency when the camera revisits regions not covered by the small window w.

3.2 TOWARDS SCENE-CONSISTENT CAMERA-CONTROLLABLE VIDEO GENERATION

The key challenge of scene-consistent camera-controllable video generation lies in reconciling two competing requirements: maintaining consistency with potentially distant frames that share spatial proximity (when the camera returns to similar viewpoints), while evolving dynamic content naturally from the recent temporal context. Ideally, conditioning on *all* frames V_{in} would ensure optimal global spatial consistency. However, this quickly becomes impractical as the sequence grows, since standard self-attention incurs quadratic time/memory in the sequence length.

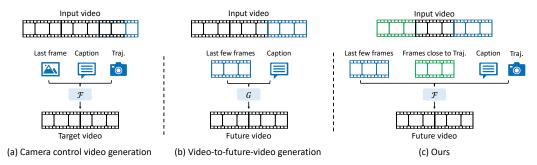


Figure 1: Comparison of existing architectures. (a) Camera-controllable methods condition on a single frame and camera trajectory. (b) Video-to-future-video methods use the last few frames of the input video when generating the future video for temporal continuity, but fail to maintain long-term spatial consistency when revisiting viewpoints unseen in the given few frames. (c) Our approach combines temporal conditioning (last few frames) with spatial conditioning (spatially adjacent frames) to achieve scene-consistent generation with precise camera control.

Dual spatio-temporal sliding window strategy. Instead of increasing the temporal window size w of the existing video-to-future-video generation methods, we introduce a dual sliding window strategy that conditions on frames selected along both temporal and spatial axes (Fig. 1-(c)). Beyond the standard temporal window that captures recent motion dynamics, we add a spatial window that retrieves frames sharing similar 3D viewpoints, regardless of their temporal distance:

$$\mathbf{V}_{\text{out}} = \mathcal{F}(\tilde{\mathbf{V}}_{\text{in}}, \mathcal{T}, \mathbf{C}), \text{ where } \tilde{\mathbf{V}}_{\text{in}} = \{\text{Temporal}(w)\} \cup \{\text{Spatial}(T)\},$$
 (4)

where the model \mathcal{F} generates a future sequence \mathbf{V}_{out} conditioned on Temporal(w), last w frames of the input video $\mathbf{V}_{\text{in}}[L-w:L]$, and Spatial(T), the T retrieved frames from the entire input sequence based on viewpoint similarity to the target viewpoint \mathbf{C} . This dual conditioning enables the model to reference distant frames that observe the same spatial regions, maintaining scene consistency without processing all L input frames.

While this dual conditioning is conceptually appealing, naïvely retrieving and providing spatially adjacent frames directly would be problematic for our task. Since we aim to generate future content beyond the input's temporal boundary, directly conditioning on frames from earlier timestamps would incorrectly preserve dynamic elements (e.g., a walking person from frame 50 should not necessarily reappear at that same location when generating frame 200). The spatial conditioning must therefore provide only the persistent scene structure while excluding dynamic content. Rather than retrieving individual frames, we introduce a **3D scene memory** $\mathcal M$ that represents exclusively the *static* geometry extracted from all spatially relevant frames.

3.3 3D Scene Memory Construction

Our 3D scene memory must efficiently encode spatial relationships across all L frames while extracting only persistent static geometry. To construct the 3D scene memory, we leverage dynamic SLAM frameworks (Li et al., 2024; Zhang et al., 2024) to estimate camera poses and reconstruct 3D structure:

$$(\hat{\mathbf{C}}, \mathbf{P}) = \mathcal{D}_{\mathsf{SLAM}}(V_{\mathsf{in}}),\tag{5}$$

where $\hat{\mathbf{C}} = \{\hat{C}_i\}_{i=1}^L$ are the estimated camera poses, \mathbf{P} represents the aggregated 3D point cloud from the L input frames, and $\mathcal{D}_{\text{SLAM}}(\cdot)$ represents the dynamic SLAM framework. This SLAM integration is effective in that it not only estimates the camera parameters of the input frames but also reconstructs the 3D structure of the scene, which can be further utilized to represent the 3D static geometry.

While the camera poses $\hat{\mathbf{C}}$ enable efficient spatial retrieval by comparing viewpoint similarity with the target trajectory \mathbf{C} , the aggregated 3D point cloud \mathbf{P} still contains both static and dynamic regions. Thus, we now explain our full pipeline on how to identify dynamic regions and only maintain the persistent static geometry of the input video.

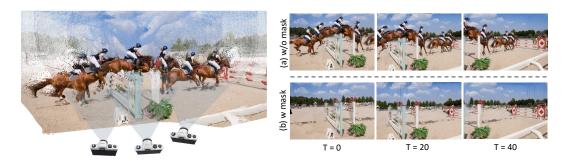


Figure 2: **Illustration of dynamic masking for static scene extraction.** When aggregating 3D points across frames, moving objects create ghosting artifacts if not properly masked. (a) Without masking, dynamic elements (horses and riders) appear frozen at multiple positions, severely degrading the warped views. (b) With our dynamic masking pipeline, these elements are identified and excluded, resulting in clean static-only point clouds that can be reliably warped to new viewpoints.

Dynamic masking for static scene extraction. Naïvely aggregating points across frames creates ghosting artifacts where moving objects appear frozen at multiple positions, as shown in Fig. 2-(a). We address this through a comprehensive three-stage masking pipeline that identifies and excludes all dynamic content as depicted in Fig. 3.

We begin with pixel-level motion detection following MonST3R (Zhang et al., 2024). For each frame pair, we compute optical flow using SEA-RAFT (Wang et al., 2024a) (Flow_{optical}) and compare it against the flow induced by camera motion alone (Flow_{warp}). Regions where the L1 difference exceeds a specific threshold τ are marked as potentially dynamic:

$$M_i^{\text{pixel}} = 1 \left[\| \text{Flow}_{\text{optical}} - \text{Flow}_{\text{warp}} \|_1 > \tau \right]. \tag{6}$$

However, pixel-level detection captures motion only at specific instants and misses complete object boundaries. We therefore propagate these sparse detections to full objects using SAM2 (Ravi et al., 2024), where we sample points from dynamic pixels in the first frame for prompts. Yet this approach still has limitations: static objects that begin moving in later frames may not be detected if they appear static initially.

Our solution employs backward tracking with CoTracker3 (Karaev et al., 2024) to aggregate motion evidence across the entire sequence. From the sampled points in each frame obtained from our pixel-level motion detection, we track these points from all frames back to t=0, capturing motions of objects that move at any point. These aggregated points are used to prompt the final SAM2 pass, producing complete object-level masks $M_i^{\rm obj}$ that cleanly separate all dynamic content (Fig. 2-(b)). With the full dynamic mask, we can now obtain the static-only 3D geometry ${\bf P}_{\rm static}$:

$$\mathbf{P}_{\text{static}} = \bigcup_{i=1}^{L} \mathbf{P}_i \odot (1 - M_i^{\text{obj}}). \tag{7}$$

From the constructed static-only 3D geometry $\mathbf{P}_{\text{static}}$ with our proposed dynamic masking strategy, we now obtain the 3D scene memory:

$$\mathcal{M} = (\hat{\mathbf{C}}, \mathbf{P}_{\text{static}}), \tag{8}$$

where we now explain how this 3D scene memory \mathcal{M} can be used for scene-consistent camera-controllable video generation in the following section.

3.4 3D Scene Prompting

Having constructed the static-only 3D representation $\mathbf{P}_{\text{static}}$, rather than naïvely retrieving T frames from the input video based on viewpoint similarity, we synthesize static-only spatial frames through the projection of $\mathbf{P}_{\text{static}}$. For each target camera pose $C_t \in \mathbf{C}$, we generate the corresponding spatial frame by projecting the static points from the most spatially relevant input frames:

$$Spatial(t) = \Pi(K \cdot C_t \cdot \mathbf{P}_{static}^{(n)}), \tag{9}$$

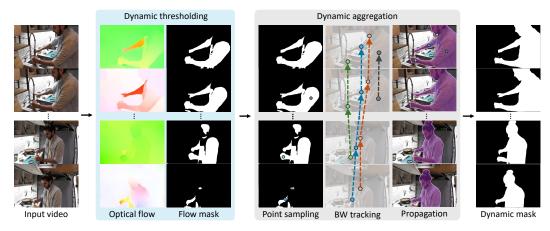


Figure 3: **Dynamic masking.** A three-stage pipeline refines dynamic region detection to produce complete object-level masks: (1) optical-flow differences detect pixel-level motion (Dynamic thresholding); (2) sample points from these regions for all frames and perform backward tracking (BW tracking) with CoTracker3 (Karaev et al., 2024) to aggregate motion evidence across all frames back to t=0 (dynamic aggregation), capturing objects that move at any time; (3) propagate aggregated points in the first frame to the entire video using SAM2 (Ravi et al., 2024). The resulting dynamic masks cleanly separate moving elements (people, objects) from the static background, enabling construction of the static-only point cloud $\mathbf{P}_{\text{static}}$.

where $\mathbf{P}_{\text{static}}^{(n)} \subset \mathbf{P}_{\text{static}}$ contains points from the top-n spatially adjacent input frames to C_t , $\Pi(\cdot)$ denotes perspective projection, and K is the camera intrinsic matrix. The complete spatial conditioning becomes $\text{Spatial}(T) = \{\text{Spatial}(t)\}_{t=1}^T \in \mathbb{R}^{T \times H \times W \times 3}$, where spatial adjacency is calculated by field-of-view overlap.

This projection-based approach ensures only static content appears in conditioning while providing geometrically consistent views aligned to target poses. Notably, the static point cloud aggregates information from multiple viewpoints, potentially filling regions occluded by dynamic objects. These projected views serve as 3D scene prompts that provide explicit guidance about persistent scene structure, enabling precise camera control without additional encoding modules.

The projected views $\operatorname{Spatial}(T)$ serve as what we term $\operatorname{3D}$ scene prompts—they provide the model with explicit guidance about the persistent scene structure. By conditioning on both $\operatorname{Temporal}(w)$ and $\operatorname{Spatial}(T)$, our framework effectively enables scene-consistent camera-controllable video generation with computational efficiency while preserving the prior for high-quality video synthesis.

4 EXPERIMENTS

4.1 IMPLEMENTATION DETAILS

Model architecture. We build upon CogVideoX-I2V-5B (Yang et al., 2024), extending its single-image conditioning to accept dual spatio-temporal inputs with minimal architectural changes. The key modification is repurposing the existing image conditioning channel to accept concatenated latents from both temporal frames and spatial projections. Specifically, we provide the last w=9 frames from $\mathbf{V}_{\rm in}$ as temporal conditioning and T projected views from the static point cloud as spatial conditioning. This enables the DiT backbone to remain entirely unchanged, preserving all pretrained video priors. Both conditions are encoded through the frozen 3D VAE and concatenated channel-wise such that $\mathbf{Z}_{\rm cond} = {\rm Concat}[\mathcal{E}({\rm Temporal}(w)), \mathcal{E}({\rm Spatial}(T))]$.

Fine-tuning. We fully fine-tune the model for a total of 4K iterations with a batch size of 8 using 4 H100 GPUs, which required approximately 48 hours. We used the 16-bit Adam optimizer with a learning rate of 1×10^{-5} , and adopted the same hyperparameter settings as those used in the training of CogVideoX (Yang et al., 2024). For the temporal sliding window, we provide the last 9 frames of the input video, setting w = 9. For the projection of top-n spatially adjacent views, we set n = 7.

Methods	RealEstate10K				DynPose-100K			
Methous	PSNR↑	SSIM↑	LPIPS↓	MEt3R↓	PSNR↑	SSIM↑	LPIPS↓	MEt3R↓
DFoT Song et al. (2025)	18.3044	0.5960	0.3077	0.181164	12.1471	0.3040	0.4172	0.183202
3DScenePrompt (Ours)	20.8932	0.7171	0.2120	0.040843	13.0468	0.3666	0.3812	0.124189

Table 1: **Evaluation of spatial and geometric consistency.** We compare DFoT and Ours on the RealEstate10K (Zhou et al., 2018) and DynPose-100K (Rockwell et al., 2025) datasets. For spatial consistency, we evaluate PSNR, SSIM, and LPIPS on revisited camera trajectories, while for geometric consistency, we report the MEt3R (Asim et al., 2025) metric.

Experimental settings. We evaluate our method across four key aspects: camera controllability, video quality, scene consistency, and geometric consistency. Since no prior work directly addresses scene-consistent camera-controllable video generation, we compare against two categories of baselines: (1) camera-controllable methods (CameraCtrl (He et al., 2024), MotionCtrl (Wang et al., 2024b), FloVD (Jin et al., 2025), AC3D (Bahmani et al., 2025)) for camera control and video quality metrics, and (2) DFoT (Song et al., 2025), which attempts scene-consistent camera-controllable generation, for spatial and geometric consistency metrics.

We primarily evaluate on 1,000 dynamic videos from DynPose-100K (Rockwell et al., 2025). For scene consistency evaluation, we additionally test on 1,000 static videos from RealEstate10K (Zhou et al., 2018), as static scenes provide clearer spatial consistency assessment.

4.2 Scene-consistent video generation

Evaluation Protocol. As mentioned in Section 3.1, one of the unique and key challenges in scene-consistent camera-controllable video generation is maintaining spatial consistency over extended durations. From a given input video, we evaluate spatial consistency by generating camera trajectories that revisit the viewpoints in the given video. By matching frames in the generated video and the input video that share the same viewpoint, we calculate PSNR, SSIM, and LPIPS. For RealEstate10K, we evaluate the whole image, whereas we only evaluate the static regions by masking out the dynamic regions for DynPose-100K. We also assess geometric consistency using Met3R (Asim et al., 2025), which measures multi-view alignment of generated frames under the recovered camera pose.

Results. As shown in Tab. 1, **3DScenePrompt** significantly outperforms DFoT across all metrics for both static and dynamic scenes. Most notably, our MEt3R evaluation error drops 77% (0.041 vs 0.181), demonstrating superior multi-view geometric alignment. While DFoT similarly tackles scene-consistent camera-controllable video generation through history guidance, their approach fails to maintain scene-consistency for long sequences due to memory constraints. In contrast, our dual spatio-temporal conditioning enables long-term scene-consistency without causing significant computational overhead. The qualitative comparisons shown in Fig. 4 also validate the effectiveness of our approach over DFoT.

4.3 CAMERA-CONTROLLABLE VIDEO GENERATION

Evaluation Protocol. We employ the evaluation protocol of previous methods (He et al., 2024; Zheng et al., 2024; Jin et al., 2025) for the camera controllability. We provide an input image along with associated camera parameters for I2V models (He et al., 2024; Wang et al., 2024b; Jin et al., 2025) and solely provide camera parameters for the T2V model (Bahmani

Table 2: Camera controllability evaluation.

Methods	DynPose-100K					
Methods	mRotErr (°)↓	mTransErr↓	mCamMC↓			
MotionCtrl Wang et al. (2024b)	3.5654	7.8231	9.7834			
CameraCtrl He et al. (2024)	3.3273	9.5989	11.2122			
FloVD Jin et al. (2025)	3.4811	11.0302	12.6202			
AC3D Bahmani et al. (2025)	3.0675	9.7044	11.1634			
DFoT Song et al. (2025)	2.3977	8.0866	9.2330			
3DScenePrompt (Ours)	2.3772	7.4174	8.6352			

et al., 2025). For our model, we provide the last 9 frames of the input video together with the camera parameters. To evaluate how faithfully the generated video follows the camera condition, we estimate camera parameters from the synthesized video using MegaSAM (Li et al., 2024), and compare the estimated camera parameters against the condition camera trajectory C.

Figure 4: Visualization of generated videos following trajectories that revisit early frames in the input video. We visualize and compare frames obtained from DFoT (Song et al., 2025) and Ours. We condition both DFoT and ours to generate a frame of a viewpoint that aligns with the viewpoint within the input. The comparison shows that ours shows much more consistent generation, whereas DFoT fails to generate scene-consistent frames mainly due to the limited number of frames it can condition on.

Table 3: **Evaluation of video generation quality.** We assess the quality of generated videos using FVD and VBench++ scores. For FVD, lower values indicate higher video quality. For VBench++ scores, higher values indicate better performance for all metrics. All VBench++ scores are normalized.

Methods	DynPose-100K								
Wedious	FVD	Overall	Subject	Bg	Aesthetic	Imaging	Temporal	Motion	Dynamic
		Score	Consist	Consist	Quality	Quality	Flicker	Smooth	Degree
MotionCtrl Wang et al. (2024b)	1017.4247	0.5625	0.5158	0.7093	0.3157	0.3149	0.8297	0.8432	0.7900
CameraCtrl He et al. (2024)	737.0506	0.6280	0.6775	0.8238	0.3736	0.3888	0.6837	0.6955	0.9900
FloVD Jin et al. (2025)	171.2697	0.7273	0.7964	0.8457	0.4722	0.5546	0.7842	0.8364	0.9900
AC3D Bahmani et al. (2025)	281.2140	0.7428	0.8360	0.8674	0.4766	0.5381	0.8020	0.8673	1.0000
3DScenePrompt (Ours)	127.4758	0.7747	0.8669	0.8727	0.4990	0.5964	0.8551	0.9260	1.0000

The comparison between the estimated and input camera parameters is quantified using three metrics: mean rotation error (mRotErr), mean translation error (mTransErr), and mean error in the camera extrinsic matrices (mCamMC). For the generated video, we also assess video synthesis performance using the Fréchet Video Distance (FVD)Skorokhodov et al. (2022) and seven metrics from VBench++Huang et al. (2024): subject consistency, background consistency, aesthetic quality, imaging quality, temporal flickering, motion smoothness, and dynamic degree.

Results. We first evaluate camera controllability and compare our method with competitive baselines. As shown in Tab. 2, our approach consistently outperforms existing methods, indicating **3DScenePrompt** is capable of generating videos with precise camera control. We then assess the overall video quality (Tab. 3) and provide qualitative comparisons (Fig. 5). As observed in Tab. 3, our method achieves the best generation quality across all metrics for dynamic video generation, which is further supported by the visual results in Fig. 5.

4.4 ABLATION STUDIES

We analyze two critical components of our framework: the dynamic masking strategy that separates static and dynamic elements, and the number of spatially adjacent frames n retrieved for spatial conditioning. Tab. 4 demonstrates

Table 4: Ablation study.

Methods	Dynamic	DynPose-100K						
Methods	mask M	PSNR↑	SSIM↑	LPIPS↓	MEt3R↓			
Ours $(n=4)$	/	13.0382	0.3733	0.3758	0.124893			
Ours $(n = L)$	/	13.0206	0.3631	0.3810	0.123507			
Ours $(n=7)$	Х	12.2304	0.3063	0.3821	0.134885			
Ours $(n=7)$	/	13.0468	0.3666	0.3812	0.124189			



Figure 5: **Visualization of scene-consistent camera-controllable video generation.** Comparison of different methods for generating videos from the same input (shown in Input Video) that follow the camera trajectory shown in GT, which is the ground truth future video. Our method best preserves scene consistency with the input video. Note the red-boxed regions in the left scene: while the input video shows a white wall, competing methods either lose scene detail or fail to maintain the original scene structure. In contrast, our approach accurately remembers the white wall and maintains consistent scene elements throughout generation. In addition, when compared with the GT Future Video, ours best follows the camera condition, effectively verifying the strength of our framework for scene-consistent camera-controllable video generation.

the impact of varying n and the necessity of dynamic masking. Without dynamic masking (3rd row), the model suffers significantly across all, showing a large drop of PSNR of approximately 0.8dB and also an increase of MEt3R error. This degradation occurs because unmasked dynamic objects create ghosting artifacts when warped to new viewpoints, corrupting the spatial conditioning. Regarding the number of spatially adjacent frames, we find that performance stabilizes around n=7, with minimal improvements beyond this point. This suggests that 7 frames provide sufficient spatial context while maintaining computational efficiency.

5 Conclusion

In this work, we introduced **3DScenePrompt**, a framework for scene-consistent camera-controllable video generation. By combining dual spatio-temporal conditioning with a static-only 3D scene memory constructed through dynamic SLAM and our dynamic masking strategy, we enable generating continuations from arbitrary-length videos while preserving scene geometry and allowing natural motion evolution. Extensive experiments demonstrate superior performance in camera controllability, scene consistency, and generation quality compared to existing methods. Our approach opens new possibilities for long-form video synthesis applications where maintaining both spatial consistency and precise camera control is essential.

REFERENCES

- Niket Agarwal, Arslan Ali, Maciej Bala, Yogesh Balaji, Erik Barker, Tiffany Cai, Prithvijit Chattopadhyay, Yongxin Chen, Yin Cui, Yifan Ding, et al. Cosmos world foundation model platform for physical ai. *arXiv preprint arXiv:2501.03575*, 2025.
- Mohammad Asim, Christopher Wewer, Thomas Wimmer, Bernt Schiele, and Jan Eric Lenssen. Met3r: Measuring multi-view consistency in generated images. *arXiv preprint arXiv:2501.06336*, 2025.
 - Sherwin Bahmani, Ivan Skorokhodov, Aliaksandr Siarohin, Willi Menapace, Guocheng Qian, Michael Vasilkovsky, Hsin-Ying Lee, Chaoyang Wang, Jiaxu Zou, Andrea Tagliasacchi, et al. Vd3d: Taming large video diffusion transformers for 3d camera control. *arXiv preprint arXiv:2407.12781*, 2024.
 - Sherwin Bahmani, Ivan Skorokhodov, Guocheng Qian, Aliaksandr Siarohin, Willi Menapace, Andrea Tagliasacchi, David B. Lindell, and Sergey Tulyakov. Ac3d: Analyzing and improving 3d camera control in video diffusion transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 22875–22889, June 2025.
 - 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. *arXiv preprint arXiv:2503.11647*, 2025.
 - 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.
 - 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. URL https://openai.com/research/video-generation-models-as-world-simulators.
 - Yuwei Guo, Ceyuan Yang, Anyi Rao, Zhengyang Liang, Yaohui Wang, Yu Qiao, Maneesh Agrawala, Dahua Lin, and Bo Dai. Animatediff: Animate your personalized text-to-image diffusion models without specific tuning. *arXiv preprint arXiv:2307.04725*, 2023.
 - Hao He, Yinghao Xu, Yuwei Guo, Gordon Wetzstein, Bo Dai, Hongsheng Li, and Ceyuan Yang. Cameractrl: Enabling camera control for text-to-video generation. *arXiv* preprint *arXiv*:2404.02101, 2024.
 - Haoran He, Yang Zhang, Liang Lin, Zhongwen Xu, and Ling Pan. Pre-trained video generative models as world simulators. *arXiv preprint arXiv:2502.07825*, 2025.
 - Ziqi Huang, Fan Zhang, Xiaojie Xu, Yinan He, Jiashuo Yu, Ziyue Dong, Qianli Ma, Nattapol Chanpaisit, Chenyang Si, Yuming Jiang, et al. Vbench++: Comprehensive and versatile benchmark suite for video generative models. *arXiv preprint arXiv:2411.13503*, 2024.
- Wonjoon Jin, Qi Dai, Chong Luo, Seung-Hwan Baek, and Sunghyun Cho. Flovd: Optical flow meets video diffusion model for enhanced camera-controlled video synthesis. *arXiv preprint arXiv:2502.08244*, 2025.
- Nikita Karaev, Iurii Makarov, Jianyuan Wang, Natalia Neverova, Andrea Vedaldi, and Christian Rupprecht. Cotracker3: Simpler and better point tracking by pseudo-labelling real videos. *arXiv* preprint arXiv:2410.11831, 2024.
- Václav Knapp and Maty Bohacek. Synthetic human action video data generation with pose transfer.
 In Synthetic Data for Computer Vision Workshop@ CVPR 2025, 2025.
 - Zhengqi Li, Richard Tucker, Forrester Cole, Qianqian Wang, Linyi Jin, Vickie Ye, Angjoo Kanazawa, Aleksander Holynski, and Noah Snavely. Megasam: Accurate, fast, and robust structure and motion from casual dynamic videos. *arXiv preprint arXiv:2412.04463*, 2024.

- Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham
 Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2: Segment anything in images
 and videos. arXiv preprint arXiv:2408.00714, 2024.
 - 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. *arXiv* preprint arXiv:2503.03751, 2025.
 - 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 *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 12444–12455, 2025.
 - Runway. Runway. https://runwayml.com/. Accessed: 2024-11-05.
 - Junyoung Seo, Jisang Han, Jaewoo Jung, Siyoon Jin, Joungbin Lee, Takuya Narihira, Kazumi Fukuda, Takashi Shibuya, Donghoon Ahn, Shoukang Hu, et al. Vid-camedit: Video camera trajectory editing with generative rendering from estimated geometry. *arXiv preprint arXiv:2506.13697*, 2025.
 - Ivan Skorokhodov, Sergey Tulyakov, and Mohamed Elhoseiny. Stylegan-v: A continuous video generator with the price, image quality and perks of stylegan2. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 3626–3636, 2022.
 - Kiwhan Song, Boyuan Chen, Max Simchowitz, Yilun Du, Russ Tedrake, and Vincent Sitzmann. History-guided video diffusion. *arXiv preprint arXiv:2502.06764*, 2025.
 - Yihan Wang, Lahav Lipson, and Jia Deng. Sea-raft: Simple, efficient, accurate raft for optical flow. In *European Conference on Computer Vision*, pp. 36–54. Springer, 2024a.
 - Zhouxia Wang, Ziyang Yuan, Xintao Wang, Yaowei Li, Tianshui Chen, Menghan Xia, Ping Luo, and Ying Shan. Motionctrl: A unified and flexible motion controller for video generation. In *ACM SIGGRAPH 2024 Conference Papers*, pp. 1–11, 2024b.
 - 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.
 - Mark Yu, Wenbo Hu, Jinbo Xing, and Ying Shan. Trajectorycrafter: Redirecting camera trajectory for monocular videos via diffusion models. In *Proceedings of the IEEE/CVF international conference on computer vision*, 2025.
 - Shangjin Zhai, Zhichao Ye, Jialin Liu, Weijian Xie, Jiaqi Hu, Zhen Peng, Hua Xue, Danpeng Chen, Xiaomeng Wang, Lei Yang, et al. Stargen: A spatiotemporal autoregression framework with video diffusion model for scalable and controllable scene generation. *arXiv preprint arXiv:2501.05763*, 2025.
 - Junyi Zhang, Charles Herrmann, Junhwa Hur, Varun Jampani, Trevor Darrell, Forrester Cole, Deqing Sun, and Ming-Hsuan Yang. Monst3r: A simple approach for estimating geometry in the presence of motion. *arXiv preprint arXiv:2410.03825*, 2024.
 - Lymin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 3836–3847, 2023.
 - Ruihan Zhang, Borou Yu, Jiajian Min, Yetong Xin, Zheng Wei, Juncheng Nemo Shi, Mingzhen Huang, Xianghao Kong, Nix Liu Xin, Shanshan Jiang, et al. Generative ai for film creation: A survey of recent advances. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 6267–6279, 2025.
 - Zhoutong Zhang, Forrester Cole, Zhengqi Li, Michael Rubinstein, Noah Snavely, and William T Freeman. Structure and motion from casual videos. In *European Conference on Computer Vision*, pp. 20–37. Springer, 2022.

Guangcong Zheng, Teng Li, Rui Jiang, Yehao Lu, Tao Wu, and Xi Li. Cami2v: Camera-controlled image-to-video diffusion model. *arXiv preprint arXiv:2410.15957*, 2024.

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

APPENDIX

Symbol	Meaning				
$V_{\text{in}} \in \mathbb{R}^{L \times H \times W \times 3}$	input arbitrary-length video with L frames.				
$V_{\mathrm{out}} \in \mathbb{R}^{T \times H \times W \times 3}$	generated future video with T frames.				
\mathcal{T}	text prompt for video generation for video models based on text-to-video (T2V) generation				
$\mathbf{C} = \{C_t\}_{t=1}^T$	desired camera trajectory the generated video $V_{ m out}$ should follow.				
C_t	camera extrinsics parameter where $C_t \in \mathbb{SE}(3)$.				
K	camera intrinsics parameter.				
$\mathcal{F}(\cdot)$	camera-controllable video generation framework.				
$\mathcal{G}(\cdot)$	video-to-future-video generation framework.				
$\mathbf{I}_{\mathrm{ref}}$	image condition for image-to-video (I2V) generation.				
V[x:y]	indexing operation; samples frames between frame x and (y-1).				
Temporal(w)	temporally adjacent w frames for condition.				
Spatial(T)	spatially adjacent T frames for condition.				
$ ilde{V}_{ m in}$	conditioning frames for our framework, includes both $\operatorname{Temporal}(w)$ and $\operatorname{Spatial}(T)$.				
$\mathcal{D}(\cdot)$	dynamic SLAM frameworks.				
$\Pi(\cdot)$	perspective projection operator.				
P	aggregated 3D point clouds.				
$\mathbf{P}_{ ext{static}}$	aggregated 3D point clouds only from static regions.				
\mathcal{M}	3D scene memory composed of camera extrinsics C_t and static point clouds $\mathbf{P}_{\text{static}}$.				
M_i^{obj}	object-level masks representing dynamic regions of frame i .				
$\mathcal{E}(\cdot)$	3D VAE.				

ADDITIONAL EXPERIMENTAL RESULTS

Additional qualitative results are provided in the supplementary zip file.

REPRODUCIBILITY STATEMENT В

As mentioned in Section 3.1, our model builds upon the open-sourced CogvideoX-I2V-5B (Yang et al., 2024) model, where each of the processes for dynamic masking is also detailedly explained. We will also make all the codes publicly available.

USE OF LARGE LANGUAGE MODELS

In accordance with the ICLR 2026 submission policy, we disclose that we used Large Language Models to assist in grammar correction for the writing in this manuscript.