

Novelty Assessment Report

Paper: Lyra: Generative 3D Scene Reconstruction via Video Diffusion Model Self-Distillation

PDF URL: <https://openreview.net/pdf?id=tIVCfVnIH0>

Venue: ICLR 2026 Conference Submission

Year: 2026

Report Generated: 2025-12-29

Abstract

The ability to generate virtual environments is crucial for applications ranging from gaming to physical AI domains such as robotics, autonomous driving, and industrial AI. Current learning-based 3D reconstruction methods rely on the availability of captured real-world multi-view data, which is not always readily available. Recent advancements in video diffusion models have shown remarkable imagination capabilities, yet their 2D nature limits the applications to simulation where a robot needs to navigate and interact with the environment. In this paper, we propose a self-distillation framework that aims to distill the implicit 3D knowledge in the video diffusion models into an explicit 3D Gaussian Splatting (3DGS) representation, eliminating the need for multi-view training data. Specifically, we augment the typical RGB decoder with a 3DGS decoder, which is supervised by the output of the RGB decoder. In this approach, the 3DGS decoder can be purely trained with synthetic data generated by video diffusion models. At inference time, our model can synthesize 3D scenes from either a text prompt or a single image for real-time rendering. Our framework further extends to dynamic 3D scene generation from a monocular input video. Experimental results show that our framework achieves state-of-the-art performance in static and dynamic 3D scene generation. Video results: <https://anonlyra.github.io/anonlyra>

Disclaimer

This report is **AI-GENERATED** using Large Language Models and WisPaper (a scholar search engine). It analyzes academic papers' tasks and contributions against retrieved prior work. While this system identifies **POTENTIAL** overlaps and novel directions, **ITS COVERAGE IS NOT EXHAUSTIVE AND JUDGMENTS ARE APPROXIMATE**. These results are intended to assist human reviewers and **SHOULD NOT** be relied upon as a definitive verdict on novelty.

Note that some papers exist in multiple, slightly different versions (e.g., with different titles or URLs). The system may retrieve several versions of the same underlying work. The current automated pipeline does not reliably align or distinguish these cases, so human reviewers will need to disambiguate them manually.

If you have any questions, please contact: mingzhang23@m.fudan.edu.cn

Core Task Landscape

This paper addresses: **Generative 3D Scene Reconstruction from Single Images or Videos**

A total of **50 papers** were analyzed and organized into a taxonomy with **15 categories**.

Taxonomy Overview

The research landscape has been organized into the following main categories:

- **Single-Image 3D Scene Reconstruction**
- **Video-Based 3D Scene Reconstruction**
- **Generative Model-Based 3D Synthesis**
- **Novel View Synthesis and 4D Scene Generation**
- **Specialized and Application-Driven Reconstruction**
- **Methodological Foundations and Surveys**

Complete Taxonomy Tree

- Generative 3D Scene Reconstruction from Single Images or Videos Survey Taxonomy
- Single-Image 3D Scene Reconstruction
 - Object-Aware and Compositional Scene Reconstruction (4 papers)
 - [1] Object-Aware 3D Scene Reconstruction from Single 2D Images of Indoor Scenes (Mingyun Wen, 2023) [View paper](#)
 - [3] Gen3dsr: Generalizable 3d scene reconstruction via divide and conquer from a single view (Andreea Ardelean, 2025) [View paper](#)
 - [9] Depr: Depth guided single-view scene reconstruction with instance-level diffusion (Zhao Qingcheng, 2025) [View paper](#)
 - [10] Coherent 3D Scene Diffusion From a Single RGB Image (Dahnert, 2024) [View paper](#)
 - Holistic Scene Reconstruction with Generative Priors (3 papers)
 - [11] Wonderworld: Interactive 3d scene generation from a single image (Hong-Xing Yu, 2025) [View paper](#)
 - [13] A recipe for generating 3d worlds from a single image (Schwarz, 2025) [View paper](#)
 - [41] VistaDream: Sampling multiview consistent images for single-view scene reconstruction (Wang Hai-ping, 2024) [View paper](#)
 - Domain-Specific Single-Image Reconstruction (2 papers)
 - [6] CGAN-Based Forest Scene 3D Reconstruction from a Single Image (Yuan Li, 2024) [View paper](#)
 - [42] A monocular thoracoscopic 3D scene reconstruction framework based on NeRF. (Juntao Han, 2025) [View paper](#)
- Video-Based 3D Scene Reconstruction
 - Dynamic Scene Reconstruction from Monocular Video (6 papers)
 - [4] Dynamic view synthesis from dynamic monocular video (Gao Chen, 2021) [View paper](#)
 - [7] Dreamo: Articulated 3d reconstruction from a single casual video (Tao Tu, 2025) [View paper](#)
 - [8] Dreamscene4d: Dynamic multi-object scene generation from monocular videos (Wen-Hsuan Chu, 2024) [View paper](#)
 - [21] Dynamic 3D Scene Reconstruction from Classroom Videos (Sebastian Janampa, 2024) [View paper](#)
 - [28] Dynibar: Neural dynamic image-based rendering (Zhengqi Li, 2023) [View paper](#)
 - [32] Neural scene flow fields for space-time view synthesis of dynamic scenes (Zhengqi Li, 2021) [View paper](#)
 - Real-Time and Mobile Video Reconstruction (2 papers)
 - [27] Generating Multi-View Action Data from a Monocular Camera Video by Fusing Human Mesh Recovery and 3D Scene Reconstruction (Hyunsu Kim, 2025) [View paper](#)
 - [38] Mobile3DRecon: Real-time monocular 3D reconstruction on a mobile phone (Xingbin Yang, 2020) [View paper](#)
 - Specialized Object and Human Reconstruction from Video (2 papers)
 - [19] High-fidelity facial avatar reconstruction from monocular video with generative priors (Yunpeng Bai, 2023) [View paper](#)

- [26] Reconstructing Hand-Held Objects in 3D from Images and Videos (Wu, 2024) [View paper](#)
- Generative Model-Based 3D Synthesis
 - 3D-Aware Generative Models for Objects and Faces (7 papers)
 - [16] Dual Encoder GAN Inversion for High-Fidelity 3D Head Reconstruction from Single Images (Bilecen, 2024) [View paper](#)
 - [18] Singraf: Learning a 3d generative radiance field for a single scene (Minjung Son, 2023) [View paper](#)
 - [25] Learning single-image 3d reconstruction by generative modelling of shape, pose and shading (Henderson, 2020) [View paper](#)
 - [29] Progressive learning of 3d reconstruction network from 2d gan data (AyÅ¼egÅ¼l DÅ¼andar, 2023) [View paper](#)
 - [31] Ganfit: Generative adversarial network fitting for high fidelity 3d face reconstruction (BarÅ¼ GeÅ¼er, 2019) [View paper](#)
 - [45] G3DR: Generative 3D Reconstruction in ImageNet (Pradyumna Reddy, 2024) [View paper](#)
 - [49] Disentangled3d: Learning a 3d generative model with disentangled geometry and appearance from monocular images (Ayush Tewari, 2022) [View paper](#)
 - Scene-Level Generative Models with Diffusion Priors (3 papers)
 - [12] Video Perception Models for 3D Scene Synthesis (Huang Rui, 2025) [View paper](#)
 - [43] 3inGAN: Learning a 3D generative model from images of a self-similar scene (Animesh Karnewar, 2022) [View paper](#)
 - [48] Sampling 3d gaussian scenes in seconds with latent diffusion models (Henderson, 2024) [View paper](#)
- Novel View Synthesis and 4D Scene Generation
 - Sparse-View and Single-View Novel View Synthesis (5 papers)
 - [2] Generative Models for View Synthesis From Sparse Images (Trevithick, 2025) [View paper](#)
 - [5] Single view generalizable 3D reconstruction based on 3D Gaussian splatting (Kun Fang, 2025) [View paper](#)
 - [15] Spatialcrafter: Unleashing the imagination of video diffusion models for scene reconstruction from limited observations (Zhang, 2025) [View paper](#)
 - [17] 5d light field synthesis from a monocular video (K. Bae, 2021) [View paper](#)
 - [24] Single-shot scene reconstruction (Sergey Zakharov, 2021) [View paper](#)
 - 4D Scene Generation from Video Diffusion Models ★ (6 papers)
 - [0] Lyra: Generative 3D Scene Reconstruction via Video Diffusion Model Self-Distillation (Anon et al., 2026) [View paper](#)
 - [34] 4dnex: Feed-forward 4d generative modeling made easy (Chen Zhaoxi, 2025) [View paper](#)
 - [36] Geo4d: Leveraging video generators for geometric 4d scene reconstruction (Jiang, 2025) [View paper](#)
 - [37] Videoscene: Distilling video diffusion model to generate 3d scenes in one step (Hanyang Wang, 2025) [View paper](#)
 - [44] 4real: Towards photorealistic 4d scene generation via video diffusion models (Yu Heng, 2024) [View paper](#)
 - [47] Holotime: Taming video diffusion models for panoramic 4d scene generation (Haiyang Zhou, 2025) [View paper](#)
 - Controllable and Interactive Scene Synthesis (4 papers)
 - [14] Generative camera dolly: Extreme monocular dynamic novel view synthesis (Van Hoorick, 2024) [View paper](#)
 - [33] Collaborative video diffusion: Consistent multi-video generation with camera control (Kuang, 2024) [View paper](#)
 - [39] Wonderverse: Extendable 3d scene generation with video generative models (Feng Hao, 2025) [View paper](#)
 - [40] FluidNexus: 3D fluid reconstruction and prediction from a single video (Yue Gao, 2025) [View paper](#)
- Specialized and Application-Driven Reconstruction
 - Autonomous Driving Scene Reconstruction (2 papers)
 - [22] Dreamdrive: Generative 4d scene modeling from street view images (Jiageng Mao, 2025) [View paper](#)
 - [35] DriveGen3D: Boosting Feed-Forward Driving Scene Generation with Efficient Video Diffusion (Wang Wei-jie, 2025) [View paper](#)
 - Text-to-3D Scene Generation (2 papers)
 - [20] VIST3A: Text-to-3D by Stitching a Multi-view Reconstruction Network to a Video Generator (Go, 2025) [View paper](#)
 - [23] Scene123: One prompt to 3D scene generation via video-assisted and consistency-enhanced MAE (Yiyang Yang, 2025) [View paper](#)
 - Unposed and Weakly-Supervised Reconstruction (2 papers)
 - [30] UVRM: A Scalable 3D Reconstruction Model from Unposed Videos (Kao, 2025) [View paper](#)
 - [50] Video supervised for 3D reconstruction from single image (Yijie Zhong, 2022) [View paper](#)
- Methodological Foundations and Surveys (1 papers)
 - [46] Recent Trends in 3D Reconstruction of General Non-Rigid Scenes (Raza Yunus, 2024) [View paper](#)

Narrative

Core task: Generative 3D scene reconstruction from single images or videos. The field divides into several complementary branches that reflect different input modalities and reconstruction goals. Single-Image 3D Scene Reconstruction focuses on inferring complete geometry and appearance from a single view, often leveraging learned priors to hallucinate occluded regions. Video-Based 3D Scene Reconstruction exploits temporal cues and multi-view consistency across frames to build richer representations. Generative Model-Based 3D Synthesis emphasizes learning-driven approaches that can synthesize plausible scenes from minimal input, while Novel View Synthesis and 4D Scene Generation extends these ideas to produce dynamic content and novel camera trajectories. Specialized and Application-Driven Reconstruction targets domain-specific challenges such as autonomous driving or medical imaging, and Methodological Foundations and Surveys provide theoretical underpinnings and comparative analyses. Representative works like Gen3dsr[3] and Gaussian Splatting Reconstruction[5] illustrate how different branches balance geometric fidelity with generative flexibility.

Recent activity has concentrated on bridging static and dynamic reconstruction, with many studies exploring how video diffusion models can generate temporally coherent 4D scenes. Trade-offs between geometric accuracy and visual plausibility remain central: some methods prioritize photorealistic synthesis at the cost of precise depth, while others enforce stricter geometric constraints. Within this landscape, Lyra[0] sits in the 4D Scene Generation from Video Diffusion Models cluster, alongside works such as 4dnex[34], Geo4d[36], and Videoscene[37]. Compared to 4real[44] and Holotime[47], which also tackle dynamic content, Lyra[0] emphasizes leveraging diffusion priors to generate novel viewpoints and temporal evolution from video input. This positioning reflects a broader trend toward integrating generative models with explicit scene representations, balancing the need for high-quality synthesis with the demand for controllable, geometrically consistent outputs.

Related Works in Same Category

The following **3 sibling papers** share the same taxonomy leaf node with the original paper:

1. 4dnex: Feed-forward 4d generative modeling made easy

Authors: Chen Zhaoxi, Liu Tianqi, Zhaoxi Chen, Zhuo Long, Tianqi Liu, et al. (20 authors total) | **Year/Venue:** 2025 | **URL:** [View paper](#)

Abstract

We present 4DNeX, the first feed-forward framework for generating 4D (i.e., dynamic 3D) scene representations from a single image. In contrast to existing methods that rely on computationally intensive optimization or require multi-frame video inputs, 4DNeX enables efficient, end-to-end image-to-4D generation by fine-tuning a pretrained video diffusion model. Specifically, 1) to alleviate the scarcity of 4D data, we construct 4DNeX-10M, a large-scale dataset with high-quality 4D annotations gene...

Relationship Analysis

Both papers belong to the 4D Scene Generation from Video Diffusion Models category, focusing on generating dynamic 3D scenes over time by leveraging video diffusion models. They overlap in their core approach of using video diffusion models to generate time-varying 3D representations from single images, with both employing self-distillation or fine-tuning strategies to extract 3D knowledge from 2D video models. The key difference is that Lyra uses a self-distillation framework with separate RGB and 3DGS decoders operating in latent space to generate 3D Gaussian Splatting representations, while 4DNeX fine-tunes a video diffusion model to directly generate unified 6D videos (RGB+XYZ sequences) representing dynamic point clouds, using width-wise fusion and modality-aware encoding strategies.

2. Geo4d: Leveraging video generators for geometric 4d scene reconstruction

Authors: Jiang, Zeren, Zheng, Chuanxia, Zeren Jiang, et al. (15 authors total) | **Year/Venue:** 2025 | **URL:** [View paper](#)

Abstract

We introduce Geo4D, a method to repurpose video diffusion models for monocular 3D reconstruction of dynamic scenes. By leveraging the strong dynamic priors captured by large-scale pre-trained video models, Geo4D can be trained using only synthetic data while generalizing well to real data in a zero-shot manner. Geo4D predicts several complementary geometric modalities, namely point, disparity, and ray maps. We propose a new multi-modal alignment algorithm to align and fuse these modalities, as w...

Relationship Analysis

Both papers belong to the same taxonomy category of 4D scene generation from video diffusion models, leveraging pre-trained video diffusion models to generate dynamic 3D scenes over time. They overlap in their core approach of distilling knowledge from video diffusion models into explicit 3D representations (3D Gaussian Splatting) for dynamic scene reconstruction from monocular inputs. The key difference is that Lyra focuses on a self-distillation framework where a 3DGS decoder is supervised by the RGB decoder of the same video model without requiring real-world multi-view data, while Geo4D predicts multiple geometric modalities (point maps, disparity maps, ray maps) and fuses them through multi-modal alignment optimization to achieve robust 4D reconstruction.

3. 4real: Towards photorealistic 4d scene generation via video diffusion models

Authors: Yu Heng, Wang Chaoyang, Heng Yu, Zhuang, Peiye, et al. (24 authors total) | **Year/Venue:** 2024 | **URL:** [View paper](#)

Abstract

Existing dynamic scene generation methods mostly rely on distilling knowledge from pre-trained 3D generative models, which are typically fine-tuned on synthetic object datasets. As a result, the generated scenes are often object-centric and lack photorealism. To address these limitations, we introduce a novel pipeline designed for photorealistic text-to-4D scene generation, discarding the dependency on multi-view generative models and instead fully utilizing video generative models trained on di...

Relationship Analysis

Both papers belong to the 4D Scene Generation from Video Diffusion Models category, leveraging video diffusion models to generate dynamic 3D scenes over time. They overlap in using video diffusion priors and 3D Gaussian Splatting representations for temporal scene reconstruction. However, Lyra focuses on self-distillation within a video model's latent space to train a 3DGS decoder without real multi-view data, while 4Real generates a reference video first, then reconstructs canonical 3DGS with per-frame deformations and temporal deformations to handle video inconsistencies.

Contributions Analysis

This paper presents **3 main contributions**, each analyzed against relevant prior work:

Contribution 1: Self-distillation framework for 3D scene reconstruction without multi-view data

Description: The authors propose a teacher-student framework where a camera-controlled video diffusion model (teacher) supervises a 3D Gaussian Splatting decoder (student) operating in latent space. This approach removes the requirement for real-world multi-view training datasets by generating synthetic supervision through the video model.

This contribution was assessed against **10 related papers** from the literature. Papers with potential prior art are analyzed in detail with textual evidence; others receive brief assessments.

1. RAFT-MSF: Self-supervised monocular scene flow using recurrent optimizer

URL: [View paper](#)

Brief Assessment

RAFT-MSF[74] addresses monocular scene flow estimation (3D motion fields from video), not 3D scene reconstruction or self-distillation frameworks. The candidate focuses on optical flow and depth estimation from consecutive frames, while the original paper proposes using video diffusion models to supervise 3D Gaussian Splatting decoders for static/dynamic scene generation.

2. Consistent 3d hand reconstruction in video via self-supervised learning

URL: [View paper](#)

Brief Assessment

Consistent Hand Video[75] focuses on self-supervised 3D hand reconstruction from monocular video using 2D keypoint detection supervision, not a teacher-student video diffusion framework for general 3D scene reconstruction.

3. Model-based 3d hand reconstruction via self-supervised learning

URL: [View paper](#)

Brief Assessment

Model-based Hand[73] focuses on self-supervised 3D hand reconstruction from single RGB images using 2D keypoint detection, not on self-distilling video diffusion models for general 3D scene reconstruction.

4. Reversing the cycle: self-supervised deep stereo through enhanced monocular distillation

URL: [View paper](#)

Brief Assessment

Reversing the Cycle[78] focuses on stereo depth estimation using monocular completion networks to improve traditional stereo methods, not on generating 3D Gaussian Splatting representations from video diffusion models for scene reconstruction.

5. Visual reinforcement learning with self-supervised 3d representations

URL: [View paper](#)

Brief Assessment

Visual Reinforcement Learning[77] focuses on learning 3D representations for robotic control tasks using multi-view supervision during training, not on eliminating multi-view data requirements through self-distillation of video diffusion models.

6. Self-supervised reflectance-guided 3d shape reconstruction from single-view images

URL: [View paper](#)

Brief Assessment

Reflectance-Guided Reconstruction[69] focuses on single-view 3D object reconstruction using reflectance properties and self-supervised learning, not on distilling video diffusion models for scene-level 3D generation without multi-view training data.

7. Pre-train, self-train, distill: A simple recipe for supersizing 3d reconstruction

URL: [View paper](#)

Brief Assessment

Pre-train Self-train Distill[72] focuses on learning a unified 3D reconstruction model across object categories using synthetic pre-training followed by self-training on image collections, then distilling category-specific models into a unified network. The original paper's self-distillation framework uses a video diffusion model as teacher to supervise a 3DGS decoder student in latent space for scene-level reconstruction. These are fundamentally different architectures and supervision mechanisms.

8. 3D Feature Distillation with Object-Centric Priors

URL: [View paper](#)

Brief Assessment

Feature Distillation[71] focuses on distilling 2D CLIP features into 3D representations for object-centric grounding in tabletop scenarios, not on self-distilling video diffusion models for general 3D scene generation without multi-view training data.

9. Exploiting the Potential of Self-Supervised Monocular Depth Estimation via Patch-Based Self-Distillation

URL: [View paper](#)

Brief Assessment

Patch-Based Self-Distillation[76] focuses on monocular depth estimation from single images using patch-based self-distillation for depth maps, not 3D scene reconstruction with explicit 3D representations like Gaussian Splatting. The domains and technical approaches are fundamentally different.

10. Weakly supervised monocular 3d detection with a single-view image

URL: [View paper](#)

Brief Assessment

Weakly Supervised Detection[70] focuses on monocular 3D object detection from single images using depth estimation, not 3D scene reconstruction. The candidate uses self-distillation between depth-guided and monocular detection networks for object localization, while the original paper distills video diffusion models into 3D Gaussian Splatting representations for scene generation.

Contribution 2: Extension to dynamic 4D scene generation from monocular video

Description: The method is extended to handle time-varying scenes by introducing time conditioning in the 3DGS decoder, enabling generation of dynamic 3D Gaussian representations from single-view video inputs with novel-view synthesis capabilities.

This contribution was assessed against **10 related papers** from the literature. Papers with potential prior art are analyzed in detail with textual evidence; others receive brief assessments.

1. Neural radiance flow for 4d view synthesis and video processing

URL: [View paper](#)

Brief Assessment

Neural Radiance Flow[59] focuses on learning 4D spatial-temporal representations using neural implicit representations (NeRF-based), not 3D Gaussian Splatting with time-conditioned decoders as in the original paper.

2. MonoFusion: Sparse-View 4D Reconstruction via Monocular Fusion

URL: [View paper](#)

Brief Assessment

MonoFusion[52] addresses sparse-view 4D reconstruction from multiple static cameras (4 equidistant inward-facing cameras), not monocular video generation. The original paper focuses on generating dynamic 3D scenes from a single video input using video diffusion models, while MonoFusion[52] reconstructs from synchronized multi-view captures with known camera poses.

3. Diffusion priors for dynamic view synthesis from monocular videos

URL: [View paper](#)

Brief Assessment

Diffusion Priors Dynamic[56] focuses on dynamic novel view synthesis from monocular videos using diffusion priors for supervision, not on feed-forward 4D generation from video diffusion models with time-conditioned 3DGS decoders as in the original paper.

4. HSR: holistic 3d human-scene reconstruction from monocular videos

URL: [View paper](#)

Brief Assessment

HSR[55] focuses on reconstructing static scenes with dynamic humans from monocular videos, not on generating novel dynamic 3D scenes with time-conditioned representations for novel-view synthesis as in the original paper's 4D generation framework.

5. Shape of motion: 4d reconstruction from a single video

URL: [View paper](#)

Prior Art Analysis

Shape of Motion[51] demonstrates prior work on 4D reconstruction from monocular video with novel-view synthesis capabilities. The candidate paper presents a method for reconstructing dynamic 4D scenes from single monocular videos, featuring explicit 3D motion trajectories and novel view synthesis. The paper describes representing dynamic scenes as persistent 3D gaussians with time-varying transformations, enabling both novel view synthesis and long-range 3D tracking from casually captured monocular videos. This directly challenges the novelty claim of extending methods to handle time-varying scenes from single-view video inputs, as Shape of Motion[51] already addresses this problem domain with similar capabilities.

Evidence

Evidence 1 - **Rationale:** Both papers address dynamic scene reconstruction from monocular video. Shape of Motion[51] explicitly states it reconstructs dynamic scenes from casually captured monocular videos, which directly overlaps with the original paper's claim of extending to dynamic 3D scene generation from monocular input video. - **Original:** Our framework further extends to dynamic 3d scene generation from a monocular input video. - **Candidate:** we introduce a method for reconstructing generic dynamic scenes, featuring explicit, persistent 3d motion trajectories in the world coordinate frame, from casually captured monocular videos.

Evidence 2 - **Rationale:** Both methods take video input and produce dynamic 3D representations with viewpoint control. Shape of Motion[51] recovers geometry and motion trajectories from video frames, enabling novel view synthesis, which matches the original paper's capability of interactive control in time and viewpoint. - **Original:** with a video input (bottom), lyra infers a dynamic 3dgs that offers interactive control in both time (rows) and viewpoint (columns). - **Candidate:** our method takes as input a sequence of t video frames $\{t \in \mathbb{R}^{h \times w \times 3}\}$ of a dynamic scene, the camera intrinsics $k \in \mathbb{R}^{3 \times 3}$, and world-to-camera extrinsics $e \in \mathbb{se}(3)$ of each input frame it. from these inputs, we aim to recover the geometry of the entire dynamic scene and the full-length 3d motion traj...

Evidence 3 - **Rationale:** Shape of Motion[51] presents a complete system for 4D reconstruction from monocular video using 3D gaussians with temporal motion, enabling novel view synthesis. This demonstrates that similar capabilities existed prior to the original paper's claimed extension to dynamic 4D generation. - **Original:** we further extend this self-distillation framework to dynamic 4d generation from monocular video. in this setting, the video model (teacher) provides space-time supervision, while the student learns to produce time-conditioned 3dgs representations that enable novel-view synthesis of dynamic scenes. - **Candidate:** we model the dense scene elements as a set of canonical 3d gaussians, that translate and rotate over entire video as persistent motion trajectories. we adopt explicit pointbased representation because it simultaneously allows for both (1) high-fidelity rendering in real-time and (2) full-length 3d t...

6. DrivingRecon: Large 4D Gaussian Reconstruction Model For Autonomous Driving

URL: [View paper](#)

Brief Assessment

DrivingRecon[58] focuses on autonomous driving scenarios using multi-view surround cameras, not monocular video inputs. The original paper extends video diffusion models to dynamic scenes from single-view video, while DrivingRecon[58] addresses feed-forward 4D reconstruction from temporal multi-view images in driving contexts.

7. Dreamscene4d: Dynamic multi-object scene generation from monocular videos

URL: [View paper](#)

Brief Assessment

Dreamscene4d[8] focuses on multi-object dynamic scene generation with object decomposition and tracking, while the original paper addresses time-conditioned 3DGS decoder for general dynamic scenes without object-level decomposition.

8. Feature4x: Bridging any monocular video to 4d agentic ai with versatile gaussian feature fields

URL: [View paper](#)

Brief Assessment

Feature4x[57] focuses on building interactive 4D feature fields for agentic AI tasks (segmentation, VQA, editing) from monocular video, not on generating dynamic 3D Gaussian representations with novel-view synthesis as the primary contribution. The original paper's contribution centers on self-distillation for 4D generation, while Feature4x[57] addresses feature field distillation for downstream AI tasks.

9. Vivid4D: Improving 4D Reconstruction from Monocular Video by Video Inpainting

URL: [View paper](#)

Brief Assessment

Vivid4D[54] focuses on 4D reconstruction from monocular video using video inpainting and depth-guided warping, not on generative 4D scene generation from video diffusion models with time-conditioned 3DGS decoders as in the original paper.

10. Diffuman4d: 4d consistent human view synthesis from sparse-view videos with spatio-temporal diffusion models

URL: [View paper](#)

Brief Assessment

Diffuman4d[53] focuses specifically on human performance reconstruction from sparse multi-view videos, not monocular single-view video inputs. The candidate requires multiple synchronized camera views as input, whereas the original contribution extends to dynamic scenes from a single monocular video.

Contribution 3: Latent-space 3DGS decoder for efficient multi-view processing

Description: The authors design a 3DGS decoder that operates directly in the compressed latent space of the video diffusion model rather than pixel space, enabling efficient fusion of hundreds of input views (726 frames) that would otherwise exceed GPU memory limits in existing pixel-based approaches.

This contribution was assessed against **10 related papers** from the literature. Papers with potential prior art are analyzed in detail with textual evidence; others receive brief assessments.

1. Splatflow: Multi-view rectified flow model for 3d gaussian splatting synthesis

URL: [View paper](#)

Prior Art Analysis

Splatflow[63] demonstrates that operating a 3DGS decoder in latent space to handle hundreds of input views was already proposed before the ORIGINAL paper. The candidate explicitly describes their GSDDecoder operating in latent space to process multi-view latents, and directly addresses the memory limitations of pixel-space approaches. The candidate states that pixel-space methods are 'restricted to 2-4 images at 512 x 512 resolution' or '24 images at 448 x 448 resolution', while their latent-space approach can handle k sparse views efficiently. This directly refutes the novelty claim that the ORIGINAL paper was first to design a latent-space 3DGS decoder for efficient multi-view processing.

Evidence

Evidence 1 - **Rationale:** Both papers describe feed-forward 3DGS decoders operating on latent representations from multi-view inputs, showing the concept existed in prior work. - **Original:** our feed-forward 3dgs decoder ds is designed to transform the synthesized multi-view latents generated by our video model v into an explicit 3d representation that can be rendered from arbitrary viewpoints. - **Candidate:** recently, feed-forward 3dgs methods [4, 9, 80, 81] have enabled fast 3dgs reconstruction from sparse views by training on large datasets, achieving much faster reconstruction than per-scene optimization methods [32]. leveraging this advantage, our gsdecoder $g\phi$ (parameterized by ϕ) is designed to dec...

Evidence 2 - **Rationale:** Both papers explicitly discuss operating in latent space rather than pixel space to handle multiple views efficiently, with the candidate already implementing this approach. - **Original:** previous feed-forward reconstruction frameworks generate 3d gaussians for each pixel, and thus, are limited by the number and resolution of input views they can handle. for example, gs-lrm (zhang et al., 2024e) operates on 2-4 images at 512 x 512 resolution, while anysplat (jiang et al., 2025a) is t... - **Candidate:** a straightforward approach is to design $g\phi$ to directly output the 3d gaussian parameters from the image latents obtained through the encoder e as $g\phi(\{e(i), \pi\}_{k=1}^K) = \{(\mu_j, \alpha_j, \sigma_j, c_j)\}_{j=1}^K$, where each 3d gaussian parameter includes position μ_j , opacity α_j , covariance σ_j , and color c_j in a...

Evidence 3 - **Rationale:** The candidate paper discusses the shift toward direct 3DGS generation methods that operate efficiently, establishing that latent-space approaches for handling multiple views were already being explored in the field. - **Original:** no latent-based 3dgs. previous works operating in the pixel space, such as btimer (liang et al., 2025b), are restricted to 12 input frames, while we can take up to 726 input frames. consequently, operating the 3dgs decoder in the pixel space instead of the latent space leads to out-of-memory. - **Candidate:** for 3dgs generation, several works [10, 42, 48, 82, 84, 94] leverage 2d diffusion models [26, 69, 76, 79] with score distillation sampling (sds) [66], which requires time-intensive per-scene optimization. to address this, recent studies have shifted towards direct 3dgs generation, combining diffusi...

2. latentsplat: Autoencoding variational gaussians for fast generalizable 3d reconstruction

URL: [View paper](#)

Prior Art Analysis

LatentSplat[64] demonstrates prior work on operating 3DGS decoders in compressed latent space rather than pixel space for efficient multi-view processing. The candidate paper explicitly addresses the same bottleneck identified in the original paper - that pixel-space methods cannot handle large numbers of high-resolution views due to GPU memory limits. LatentSplat[64] presents a solution by working in the latent space of a VAE, enabling processing of multiple views efficiently, which directly refutes the novelty claim that the original authors were first to design such a latent-space 3DGS decoder.

Evidence

Evidence 1 - **Rationale:** The original paper's ablation confirms that latent-space operation is essential for handling many frames. LatentSplat[64] demonstrates this capability was already established, showing the architectural innovation predates the original work. - **Original:** no latent-based 3dgs. previous works operating in the pixel space, such as btimer (liang et al., 2025b), are restricted to 12 input frames, while we can take up to 726 input frames. consequently, operating the 3dgs decoder in the pixel space instead of the latent space leads to out-of-memory. - **Candidate:** in contrast to all of the above, we provide high-quality 360o reconstructions of object-centric scenes as well as view interand extrapolation on large scenes, given only two input views.

3. FastAvatar: Instant 3D Gaussian Splatting for Faces from Single Unconstrained Poses

URL: [View paper](#)

Brief Assessment

FastAvatar[68] operates on single-image face reconstruction with a feed-forward encoder-decoder architecture, not multi-view video processing. The candidate does not address efficient fusion of hundreds of views or latent-space processing of video diffusion model outputs.

4. Langsplat: 3d language gaussian splatting

URL: [View paper](#)

Brief Assessment

Langsplat[60] focuses on 3D language field modeling using Gaussian splatting for open-vocabulary queries, not on multi-view video processing or handling hundreds of input frames. The latent-space operation in Langsplat is for compressing CLIP embeddings via a scene-specific autoencoder, not for processing multi-view video latents from a diffusion model.

5. L3dg: Latent 3d gaussian diffusion

URL: [View paper](#)

Brief Assessment

L3dg[65] focuses on unconditional generative modeling of 3D Gaussians using a VQ-VAE latent diffusion approach for scene synthesis, not on feed-forward reconstruction from multi-view video inputs with a video diffusion model teacher.

6. Styleme3d: Stylization with disentangled priors by multiple encoders on 3d gaussians

URL: [View paper](#)

Brief Assessment

Styleme3d[67] focuses on 3D Gaussian Splatting stylization using multiple encoders (DSSD, CSD, SOS, 3DG-QA) for artistic style transfer, not on latent-space decoders for multi-view reconstruction or memory-efficient processing of hundreds of frames.

7. DSplats: 3D Generation by Denoising Splats-Based Multiview Diffusion Models

URL: [View paper](#)

Brief Assessment

DSplats[66] operates in latent space for efficiency but uses a different architecture (U-Net-based diffusion model) rather than the video diffusion model framework described in the original paper. The candidate focuses on single-image 3D reconstruction via diffusion, not multi-trajectory video processing with 726 frames.

8. UniSplat: Unified Spatio-Temporal Fusion via 3D Latent Scaffolds for Dynamic Driving Scene Reconstruction

URL: [View paper](#)

Brief Assessment

UniSplat[61] operates on a 3D latent scaffold constructed from multi-view images, but this scaffold is built from geometry and semantic features in 3D voxel space, not from the compressed latent space of a video diffusion model. The original paper's decoder operates 'directly in the compressed latent space of the video diffusion model' to process 726 frames, while UniSplat[61] uses foundation models to construct 3D scaffolds and performs fusion in that 3D representation space.

9. VIST3A: Text-to-3D by Stitching a Multi-view Reconstruction Network to a Video Generator

URL: [View paper](#)

Brief Assessment

[Final Audit Failure] The model insisted on a refutation claim but failed to provide verifiable evidence after multiple retries. Marked as cannot_refute for safety. Please manually verify the candidate text.

10. Leveraging latent diffusion in 3D Gaussian splatting for novel view synthesis

URL: [View paper](#)

Brief Assessment

Latent Diffusion Splatting[62] mentions applying diffusion models in latent space versus image space, but the provided context is too limited to determine if they implement a latent-space 3DGS decoder for multi-view fusion or address the specific memory efficiency problem of processing 726 frames that the original paper solves.

Appendix: Text Similarity Detection

No high-similarity text segments were detected across any compared papers.

References

- [0] Lyra: Generative 3D Scene Reconstruction via Video Diffusion Model Self-Distillation [View paper](#)
- [1] Object-Aware 3D Scene Reconstruction from Single 2D Images of Indoor Scenes [View paper](#)
- [2] Generative Models for View Synthesis From Sparse Images [View paper](#)
- [3] Gen3dsr: Generalizable 3d scene reconstruction via divide and conquer from a single view [View paper](#)
- [4] Dynamic view synthesis from dynamic monocular video [View paper](#)
- [5] Single view generalizable 3D reconstruction based on 3D Gaussian splatting [View paper](#)
- [6] CGAN-Based Forest Scene 3D Reconstruction from a Single Image [View paper](#)
- [7] Dreamo: Articulated 3d reconstruction from a single casual video [View paper](#)
- [8] Dreamscene4d: Dynamic multi-object scene generation from monocular videos [View paper](#)
- [9] Depr: Depth guided single-view scene reconstruction with instance-level diffusion [View paper](#)
- [10] Coherent 3D Scene Diffusion From a Single RGB Image [View paper](#)
- [11] Wonderworld: Interactive 3d scene generation from a single image [View paper](#)
- [12] Video Perception Models for 3D Scene Synthesis [View paper](#)
- [13] A recipe for generating 3d worlds from a single image [View paper](#)
- [14] Generative camera dolly: Extreme monocular dynamic novel view synthesis [View paper](#)
- [15] Spatialcrafter: Unleashing the imagination of video diffusion models for scene reconstruction from limited observations [View paper](#)
- [16] Dual Encoder GAN Inversion for High-Fidelity 3D Head Reconstruction from Single Images [View paper](#)
- [17] 5d light field synthesis from a monocular video [View paper](#)
- [18] Singraf: Learning a 3d generative radiance field for a single scene [View paper](#)
- [19] High-fidelity facial avatar reconstruction from monocular video with generative priors [View paper](#)
- [20] VIST3A: Text-to-3D by Stitching a Multi-view Reconstruction Network to a Video Generator [View paper](#)
- [21] Dynamic 3D Scene Reconstruction from Classroom Videos [View paper](#)
- [22] Dreamdrive: Generative 4d scene modeling from street view images [View paper](#)
- [23] Scene123: One prompt to 3D scene generation via video-assisted and consistency-enhanced MAE [View paper](#)
- [24] Single-shot scene reconstruction [View paper](#)
- [25] Learning single-image 3d reconstruction by generative modelling of shape, pose and shading [View paper](#)
- [26] Reconstructing Hand-Held Objects in 3D from Images and Videos [View paper](#)
- [27] Generating Multi-View Action Data from a Monocular Camera Video by Fusing Human Mesh Recovery and 3D Scene Reconstruction [View paper](#)
- [28] Dynibar: Neural dynamic image-based rendering [View paper](#)
- [29] Progressive learning of 3d reconstruction network from 2d gan data [View paper](#)
- [30] UVRM: A Scalable 3D Reconstruction Model from Unposed Videos [View paper](#)
- [31] Ganfit: Generative adversarial network fitting for high fidelity 3d face reconstruction [View paper](#)
- [32] Neural scene flow fields for space-time view synthesis of dynamic scenes [View paper](#)
- [33] Collaborative video diffusion: Consistent multi-video generation with camera control [View paper](#)
- [34] 4dnex: Feed-forward 4d generative modeling made easy [View paper](#)
- [35] DriveGen3D: Boosting Feed-Forward Driving Scene Generation with Efficient Video Diffusion [View paper](#)
- [36] Geo4d: Leveraging video generators for geometric 4d scene reconstruction [View paper](#)
- [37] Videoscene: Distilling video diffusion model to generate 3d scenes in one step [View paper](#)
- [38] Mobile3DRecon: Real-time monocular 3D reconstruction on a mobile phone [View paper](#)
- [39] Wonderverso: Extendable 3d scene generation with video generative models [View paper](#)
- [40] FluidNexus: 3D fluid reconstruction and prediction from a single video [View paper](#)
- [41] VistaDream: Sampling multiview consistent images for single-view scene reconstruction [View paper](#)
- [42] A monocular thoracoscopic 3D scene reconstruction framework based on NeRF. [View paper](#)
- [43] 3inGAN: Learning a 3D generative model from images of a self-similar scene [View paper](#)
- [44] 4real: Towards photorealistic 4d scene generation via video diffusion models [View paper](#)
- [45] G3DR: Generative 3D Reconstruction in ImageNet [View paper](#)
- [46] Recent Trends in 3D Reconstruction of General Non-rigid Scenes [View paper](#)
- [47] Holotime: Taming video diffusion models for panoramic 4d scene generation [View paper](#)
- [48] Sampling 3d gaussian scenes in seconds with latent diffusion models [View paper](#)
- [49] Disentangled3d: Learning a 3d generative model with disentangled geometry and appearance from monocular images [View paper](#)
- [50] Video supervised for 3D reconstruction from single image [View paper](#)
- [51] Shape of motion: 4d reconstruction from a single video [View paper](#)
- [52] MonoFusion: Sparse-View 4D Reconstruction via Monocular Fusion [View paper](#)

- [53] Diffuman4d: 4d consistent human view synthesis from sparse-view videos with spatio-temporal diffusion models [View paper](#)
- [54] Vivid4D: Improving 4D Reconstruction from Monocular Video by Video Inpainting [View paper](#)
- [55] HSR: holistic 3d human-scene reconstruction from monocular videos [View paper](#)
- [56] Diffusion priors for dynamic view synthesis from monocular videos [View paper](#)
- [57] Feature4x: Bridging any monocular video to 4d agentic ai with versatile gaussian feature fields [View paper](#)
- [58] DrivingRecon: Large 4D Gaussian Reconstruction Model For Autonomous Driving [View paper](#)
- [59] Neural radiance flow for 4d view synthesis and video processing [View paper](#)
- [60] Langsplat: 3d language gaussian splatting [View paper](#)
- [61] UniSplat: Unified Spatio-Temporal Fusion via 3D Latent Scaffolds for Dynamic Driving Scene Reconstruction [View paper](#)
- [62] Leveraging latent diffusion in 3D Gaussian splatting for novel view synthesis [View paper](#)
- [63] Splatflow: Multi-view rectified flow model for 3d gaussian splatting synthesis [View paper](#)
- [64] latentsplat: Autoencoding variational gaussians for fast generalizable 3d reconstruction [View paper](#)
- [65] L3dg: Latent 3d gaussian diffusion [View paper](#)
- [66] DSplats: 3D Generation by Denoising Splats-Based Multiview Diffusion Models [View paper](#)
- [67] Styleme3d: Stylization with disentangled priors by multiple encoders on 3d gaussians [View paper](#)
- [68] FastAvatar: Instant 3D Gaussian Splatting for Faces from Single Unconstrained Poses [View paper](#)
- [69] Self-supervised reflectance-guided 3d shape reconstruction from single-view images [View paper](#)
- [70] Weakly supervised monocular 3d detection with a single-view image [View paper](#)
- [71] 3D Feature Distillation with Object-Centric Priors [View paper](#)
- [72] Pre-train, self-train, distill: A simple recipe for supersizing 3d reconstruction [View paper](#)
- [73] Model-based 3d hand reconstruction via self-supervised learning [View paper](#)
- [74] RAFT-MSF: Self-supervised monocular scene flow using recurrent optimizer [View paper](#)
- [75] Consistent 3d hand reconstruction in video via self-supervised learning [View paper](#)
- [76] Exploiting the Potential of Self-Supervised Monocular Depth Estimation via Patch-Based Self-Distillation [View paper](#)
- [77] Visual reinforcement learning with self-supervised 3d representations [View paper](#)
- [78] Reversing the cycle: self-supervised deep stereo through enhanced monocular distillation [View paper](#)