

LangScene-X: Reconstruct Generalizable 3D Language-Embedded Scenes with TriMap Video Diffusion

Fangfu Liu¹, Hao Li², Jiawei Chi¹, Hanyang Wang¹, Minghui Yang³, Fudong Wang³, Yueqi Duan^{1†}
¹Tsinghua University, ²NTU, ³Ant Group



Figure 1. **LangScene-X**: Given sparse views as input (e.g., as few as two images), we design a generative paradigm to build the 3D generalizable language-embedded surface fields with TriMap video diffusion and language quantized compressor (LQC), which supports open-ended language queries in any 3D scenes. For example, given a prompt like “stuffed bear” in the Teatime scene, a relevancy map can be rendered focusing on the location with maximum relevancy activation.

Abstract

Recovering 3D structures with open-vocabulary scene understanding from 2D images is a fundamental but daunting task. Recent developments have achieved this by performing per-scene optimization with embedded language information. However, they heavily rely on the calibrated dense-view reconstruction paradigm, thereby suffering from severe rendering artifacts and implausible semantic synthesis when limited views are available. In this paper, we introduce a novel generative framework, coined **LangScene-X**, to unify and generate 3D consistent multi-modality information for reconstruction and understanding. Powered by the generative capability of creating more consistent novel

observations, we can build generalizable 3D language-embedded scenes from only sparse views. Specifically, we first train a TriMap video diffusion model that can generate appearance (RGBs), geometry (normals), and semantics (segmentation maps) from sparse inputs through progressive knowledge integration. Furthermore, we propose a Language Quantized Compressor (LQC), trained on large-scale image datasets, to efficiently encode language embeddings, enabling cross-scene generalization without per-scene retraining. Finally, we reconstruct the language surface fields by aligning language information onto the surface of 3D scenes, enabling open-ended language queries. Extensive experiments on real-world data demonstrate the superiority of our LangScene-X over state-of-the-art methods in terms of quality and generalizability. Project Page:

[†]The corresponding author.

1. Introduction

As spatial learners with strong prior knowledge, humans can perceive, interpret, and understand the 3D physical world from only few-shot visual captures. Therefore, it is fundamental and crucial in computer vision to learn 3D structures from images with scene understanding, which allows to interact and query the 3D worlds through open-ended language [2, 10, 17, 48], offering a wide range of applications [36, 39, 53] such as robotics, autonomous navigation, and VR/AR.

With the rapid advancements in NeRF [31] and Gaussian Splatting [15]), recent works [17, 20, 34] have incorporated the language features from the CLIP [35] into the 3D representations to build a 3D language field that supports open-vocabulary object querying. Although they can achieve promising results in per-scene optimization with calibrated dense views (usually more than 20 views) as input, they cannot generalize to unseen scenes and suffer from severe artifacts in 3D structures with implausible semantics when encountering insufficient input views [13]. Moreover, these approaches heavily rely on scene-specific auto-encoders to compress 512-dim CLIP features into lower-dim latent space for efficient rendering. Such reliance significantly increases training time and a tendency for domain overfit, which hinders fast inference and the ability to effectively scale with large datasets, limiting the range of applications in real-world scenarios.

The heart of high-quality language-embedded 3D scenes lies in the need to integrate multimodal information (*i.e.*, semantics, geometry, and appearance) into a cohesive 3D representation. While dense, calibrated views provide abundant data for this integration, their acquisition is costly and impractical for broader adoption [13]. In contrast, constructing language-embedded fields from sparse views is inherently more challenging due to the scarcity of input information, yet it is critical for expanding applicability. The primary difficulty is extracting and fusing sufficient multimodal knowledge from limited inputs to achieve coherent 3D scene reconstruction and understanding.

To address this, we propose **LangScene-X**, a novel generative paradigm to build generalizable 3D language-embedded scenes from very sparse views (*i.e.*, as few as two images). Building upon the representation learning capabilities of generative models [1, 38, 50], our key insight is to unleash the strong generative prior to unify information for reconstruction and understanding in a single video diffusion model. Specifically, we first build a TriMap video diffusion model that can generate appearance (RGB images), geometry (normal maps), and semantics (semantic maps) from sparse inputs. To ensure 3D consistency

across generated frames and to bridge domain gaps between different modalities (RGB, normals, and segmentation masks), we meticulously design a multi-task training strategy that progressively integrates the knowledge from these diverse domains. Powered by the strong generalizability of video diffusion, our TriMap video diffusion can generate hierarchical semantic maps at varying levels of granularity. To reduce the memory cost and enhance scalability for large-scale data, we propose a generalizable Language Quantized Compressor (LQC) trained on large-scale datasets, which encodes high-dimensional language features as low-dimensional discrete indices without sacrificing essential properties. This approach eliminates the need for per-scene retraining, reduces memory overhead, and enables rapid rendering of language-embedded Gaussians. Finally, we reconstruct the language-embedded surface fields by contextually aligning the generated semantics onto the surface (generated normals) of 3D scenes (generated RGBs), enabling a broad range of downstream tasks with 3D language. We summarize the contributions of the paper as follows:

- We introduce LangScene-X, a novel generative framework that constructs generalizable 3D language-embedded fields from sparse views, which unify the information of scene reconstruction and understanding in a video diffusion model.
- We build the TriMap video diffusion model through a progressive multi-task training scheme, which can generate 3D-consistent RGB images, normal maps, and semantic maps across video frames.
- We propose a generalizable language quantized compressor (LQC) for efficient feature representation, reducing memory usage and rapid rendering while maintaining essential language features properties.
- We conduct extensive experiments to verify the efficacy of our framework and show that LangScene-X outperforms existing methods for high quality and generalizability, revealing the great potential to craft 3D language fields from a generative paradigm.

2. Related Work

2.1. Gaussian Splatting

3D Gaussian Splatting (3DGS) [15] has recently made significant strides in advancing 3D scene representation. It offers high-resolution real-time rendering capabilities that surpass traditional Neural Radiance Field (NeRF) methods. This efficiency facilitates various downstream applications [17, 27, 28, 30]. For scene surface reconstruction [11], PGSR [5] employs multi-view supervision to the detailed reconstruction of 3D surfaces and meshes. In the field of scene understanding [7, 25], LangSplat [34] and LangSurf [20] introduce Language Embeddings into

Gaussian attributes to facilitate efficient text querying in 3D space. Despite its capabilities, 3DGS struggles with accurately reconstructing scenes with limited observed views [4]. Several subsequent methods [6, 22] enhance this process by incorporating geometric regularization. Moreover, other methods [8, 21] utilize feed-forward models without per-scene optimization to perform generalized 3D Gaussian representations. Recently, ViewCrafter [52] and ReconX [26] leverage video diffusion models to interpolate nearby frames given single or two views to achieve dense view prediction and scene reconstruction. In contrast, our method offers a multi-task generation pattern that generates appearance (RGBs), geometry (normals), and semantics (semantic maps) from sparse views, promoting efficient and accurate scene reconstruction and understanding.

2.2. Video Diffusion Models

Integrating additional constraints for conditional video generation has yielded remarkable results in the field of 3D scene reconstruction/generation. Much of the earlier work in video generation has concentrated on incorporating various control signals to video diffusion models. For instance, some works [3, 12] introduce camera pose embedding into U-Net based video diffusion models [47] to drive controllable video generation. Moreover, ReconX [26] and ViewCrafter [52] integrate DUST3R [45] as explicit 3D prior for 3D consistent video generation. Additionally, DimensionX [41] achieves consistent video generation with large trajectory motion based on a powerful DiT-based video generation structure [49]. However, those works are limited to single modality generations, restricting their downstream applications. Rather than injecting additional control signals, some methods [14, 54] customize video diffusion models by fine-tuning on a set of reference videos with similar patterns, such as depth maps or motions. Yet, these methods often rely on complete video rather than sparse view inputs and involve complex network architectures and training procedures. In contrast, we propose a TriMap video diffusion model, which can generate 3D consistent videos across various domains (*i.e.*, appearance, geometry, and semantics) by progressively fusing the cross-domain knowledge into the DiT-based video diffusion model [49] without sacrificing performance.

2.3. Language Embedded Fields

The basic idea of the language-embedded Gaussian representation is to inject 2D language features from CLIP [35] and SAM [18] into Gaussians and then fascinating real-time text-guided object query in novel views. However, directly encoding 512-dim CLIP features into Gaussian primitives is memory expensive. LangSplat [34] adopts Autoencoder to encode scene-wise CLIP features into low-dimension latent for Gaussian training. OpenGaussian [46] utilizes a

coarse-to-fine feature codebook to register CLIP features into a low-dimension index. Nevertheless, such methods heavily rely on scene-specific feature compressors, which is time-consuming and hinders extension or generalization to other scenes. In contrast, our method proposes a Quantized Language Compressor for compact feature representation, which maps high-dimension features into a low-dim index while maintaining essential feature properties. Moreover, they require calibrated dense views to reconstruct the scene, leading to poor performance in the scenario of sparse or limited view input. Some generalizable methods [8, 21] directly encode multi-image inputs into 3D semantic representation through large foundation models [43] for novel view generation. GP-NeRF [21] adopts Dual Transformer to aggregate RGB and semantic features from input views into an implicit representation. LSM [8] integrates learning-based MAST3R [45] for dense language Gaussian predictions from sparse unposed images. However, these methods are limited to category-specific domains, such as objects and indoor scenes, and are prone to artifacts due to their limited model representation capabilities. Our method unleashes the power of video diffusion models in a novel way to generate 3D consistent RGB, appearance prediction, and segmentation masks simultaneously, offering scalability and generalization with improved performance.

3. Method

3.1. Overview of LangScene-X

Given N sparse views (*i.e.*, as few as two images) as input, our goal is to reconstruct and understand the underlying 3D scene (*i.e.*, construct the language-embedded surface fields). In our framework LangScene-X, we first build the TriMap video diffusion model to generate 3D consistent RGB images, normal maps, and semantic maps from sparse-view input (Sec. 3.2), which provides the dense frames for reconstruction and understanding. Then we compress the high-dimensional language features into low-dimensional discrete space through a generalizable language quantized compressor (Sec. 3.2). This eliminates per-scene retraining and enables rapid rendering of Gaussians. Finally, we reconstruct the language-embedded 3D scenes with generated and compressed information (Sec. 3.4), supporting open-ended language queries in any viewpoint. Our pipeline is depicted in Figure 2.

3.2. Building the TriMap Video Diffusion

Existing works like LERF [17] and LangSplat [34] essentially require very dense multi-view inputs with accurate calibration due to inherent per-scene optimization schemes from NeRF [31] or 3D Gaussian [15]. As a result, these methods often struggle to synthesize high-quality images with semantics from insufficient views, especially

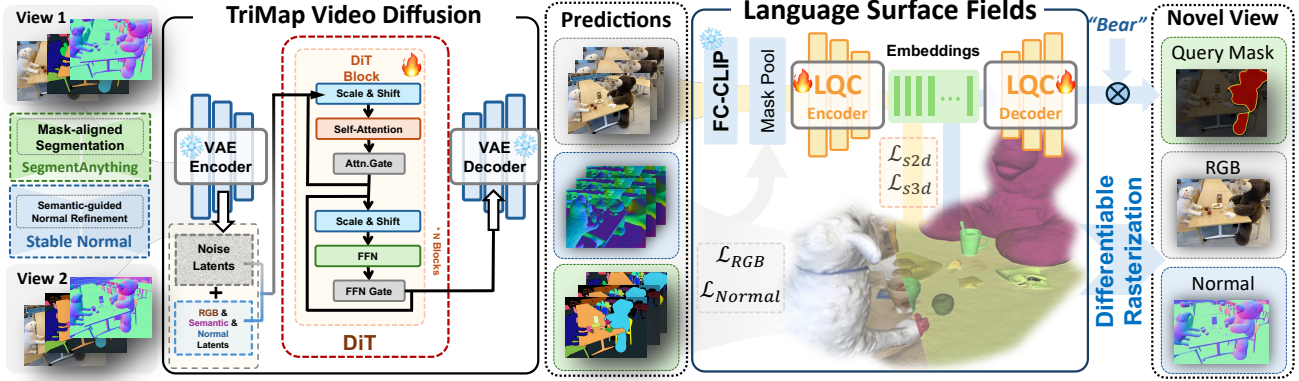


Figure 2. **Pipeline of LangScene-X.** Given two sparse-view images as input, we first generate a sequence of 3D consistent RGB images, normal maps, and segmentation maps from TriMap video diffusion model, which provides the dense frames for later 3D scene reconstruction and understanding. Then we project high-dimensional semantic features into low-dimensional discrete space through a generalizable Language Quantized Compressor (LQC). Finally, we reconstruct the 3D language-embedded scenes with generated and compressed information, which supports open-ended language queries in any viewpoint.

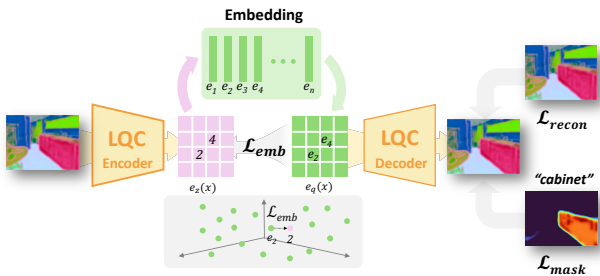


Figure 3. The illustration of Language Quantized Compressor (LQC). By leveraging learnable embedding and vector quantisation strategy, it compresses high-dimensional language features into discrete D -channel latent representation, facilitating efficient language fields reconstructions and render.

in unseen areas. Recently, video generative models have shown promise for generating frames featuring 3D structure [26, 44, 52] and perception capability [38, 40]. This inspires us to unleash the strong generative prior of large pre-trained video diffusion models to generate frames that unify the information for reconstructing (RGB images and normal maps) and understanding (segmentation masks) 3D scenes in language-embedded fields. However, it is non-trivial to achieve 3D consistency in video generation across different perception patterns [26, 40].

To address this, we build the large TriMap video diffusion model to generate RGB images, normals, and semantic maps from only sparse views (*i.e.*, as few as two images). We follow the architecture as CogVideoX [49], which is a transformer-based video diffusion model [33]. Before we train the TriMap video diffusion, we first modify the i2v architecture to support key-frame interpolation,

which is more flexible for sparse-view input. Given a T -frame video $V \in \mathbb{R}^{H \times W \times 3}$, we choose the first I_f and the last image I_e padded with zeros to obtain a condition video with the same size as V . Then, we encode the condition video with a causal VAE [49] encoder to get a latent vector concatenated with a Gaussian noise of the same size. Finally, we iteratively denoise the noise latent and apply the VAE decoder to obtain key-frame interpolation results. Keeping the architecture as unchanged as possible, we can fully leverage the pretrained prior knowledge. Let \mathcal{D}_I , \mathcal{D}_N and \mathcal{D}_S be different domain mappers \mathcal{M} for RGB images, normals and semantic maps respectively (*i.e.*, $\mathcal{M} : \mathbb{R}^{H \times W \times C} \times \mathcal{D}_\theta \rightarrow \mathbb{R}^{H \times W \times C}$) while keeping the same channels, we meticulously divide the training into four stages. (1) We first use large-scale web data to train the TriMap video diffusion to support a general key-frame interpolation task. (2) We then finetune the model with a small amount of 3D consistent video data ($\sim 10K$) to learn 3D consistency across frames. (3) Next, we annotate 200 clips of 3D consistent video data with \mathcal{D}_N to get normal videos along with RGB videos to finetune the model. (4) Finally, we apply \mathcal{D}_S to annotate 300 clips of semantic masks from 3D consistent video data and finetune with RGB and normal videos. The overall training objective is:

$$\mathcal{L}_{\text{diff}} = \mathbb{E}_{\mathbf{x}, \epsilon \sim \mathcal{N}(0, I), t, \mathcal{D}_i \in \mathcal{M}} \left[\|\epsilon - \epsilon_\theta(\mathbf{x}_t, t, \mathcal{D}_i)\|_2^2 \right], \quad (1)$$

where \mathbf{x}_t is the noise latent from the ground-truth views of the training data. Through empirical studies, such progressive knowledge integration training strategy can facilitate 3D consistency and multi-task perception in generated frames simultaneously. Moreover, powered by strong generalizability, our TriMap video diffusion can generate multi-

hierarchy masks $\{\mathbb{M}^h | h = s, m, l\}$ at inference time from only two input views segmented by \mathcal{D}_S , where s, m, l represents small, medium, and large hierarchy levels of the segmentation masks.

3.3. Language Quantized Compressor

Previous methods [20, 34] rely on scene-specific autoencoders to compress high-dimensional language features into low-dimensional features with per-scene optimization. Such a technique represents latent with continuous distribution, which struggles to represent essential language properties with extremely low-dimension compression (*i.e.*, 3-channels) on large-scale data. Meanwhile, language features are inherently discrete [42] (*i.e.*, features with same categories share consistent distributions). Therefore, discrete representations are naturally suited for compressing language features during the process of low-dimensional representation. To this end, as shown in Figure. 3, we leverage Vector Quantization into our compressor and propose a Language-Quantized Compressor (LQC) trained on large-scale datasets (COCO [23]). Specifically, we define a learnable embeddings $E = \{e_1, e_2, \dots, e_k | e_k \in \mathbb{R}^D\}$ to capture essential language properties during the medium layer of our compressor, where k denotes the size of discrete latent space and D denotes the channel of the latents. The language features x produced by off-the-shelf dense CLIP feature extractor [19] pass through the encoder to obtain $z_e(x)$, which serves as indices to calculate nearest neighbor look-up with embeddings E and map $z_e(x)$ into $1 \sim k$ embeddings:

$$z_q(x) = e_k, \text{ where } k = \operatorname{argmin}_j \|z_e(x) - e_j\|_2, \quad (2)$$

where $z_q(x)$ is the final mapped discrete latent. After that, $z_q(x)$ passes through the decoder and obtains the reconstructed feature \hat{x} . However, training such a compressor is non-trivial, as the look-up operation obstructs the gradient flow from the decoder to the encoder. To address it, we directly copy the gradient flow from decoder to encoder networks for encoder-decoder training, where :

$$\mathcal{L}_r = \|x - \operatorname{decoder}(z_e(x) + \operatorname{sg}(z_q(x) - z_e(x)))\|_2^2, \quad (3)$$

where sg means the operation of “stop gradient”. For learnable embeddings training, we utilize classic dictionary learning algorithms that push embeddings E towards encoder outputs $z_e(x)$:

$$\mathcal{L}_{emb} = \|\operatorname{sg}[z_e(x)] - E\|_2^2. \quad (4)$$

Additionally, to ensure accurate language-feature alignment, we utilize pseudo-mask supervision by applying L2 loss on text-guided activation maps between vanilla feature x and reconstructed feature \hat{x} :

$$\mathcal{L}_{mask} = \|\hat{x} \cdot T - x \cdot T\|_2^2 \quad (5)$$

The overall criterion can be summarized as follows:

$$\mathcal{L}_{lqc} = \lambda_1 \mathcal{L}_r + \lambda_2 \mathcal{L}_{emb} + \lambda_3 \mathcal{L}_{mask}, \quad (6)$$

where $\lambda_1, \lambda_2, \lambda_3$ denote loss weights. Such strategies enable high-ratio feature compression for subsequent language Gaussian training, preserving all essential language properties.

3.4. Language-Embedded Surface Fields

Building upon the 3D-consistent frames generation (RGB, normals, and semantic maps) and efficient compress pipeline, we are able to construct accurate language-embedded surface fields that facilitate efficient novel-view synthesis and open-ended 3D language-guided query. Specifically, given the generated RGB sequences $\mathbf{C} \in \mathbb{R}^{D \times H \times W \times 3}$, we initialize the sparse point clouds of the scene using DUST3R [45] and train our language surface fields using regular L2 RGB loss \mathcal{L}_{rgb} with several steps. To ensure accurate 3D language representation, we then perform joint training using geometry and semantic supervision to construct language-embedded surface fields. In practice, we leverage the powerful normal priors $\mathbf{N} \in \mathbb{R}^{D \times H \times W \times 3}$ generated by the TriMap video diffusion, we adopt a progressive normal regularization \mathcal{L}_{normal} to optimize the geometry representation of our model:

$$\mathcal{L}_{normal} = \begin{cases} \|\mathbf{N}_p - \mathbf{N}\|_1, & \text{step} < T_n \\ \|\mathbf{N}_p - \hat{\mathbf{N}}\|_1, & \end{cases}, \quad (7)$$

where \mathbf{N}_p is the rendered normal by our training fields and T_n is the hyper-parameter step. $\hat{\mathbf{N}}$ is the generated normal that filters out uncertain regions that hard to optimize ($\theta_n > \tau_{thr}$) to alleviate the impact of incorrect normal generation. The τ_{thr} is the hyper-parameter threshold and θ_n is the angle difference between \mathbf{N}_p and \mathbf{N} :

$$\theta_p = \arccos \left(\frac{\mathbf{N}_p \cdot \mathbf{N}}{\|\mathbf{N}_p\| \|\mathbf{N}\|} \right). \quad (8)$$

For semantic supervision, we utilize pre-trained dense CLIP extractor [51] to extract language features from generated images \mathbf{C} and pass through our proposed LQC to obtain 3-channel language latent. Moreover, to facilitate accurate 3D semantic representation, apart from regular L2 semantic loss, we additionally adopt a 2D and 3D clustering criterion based on the generated segmentation masks \mathbf{M} :

$$\begin{cases} \mathcal{L}_{s2d} = \|\hat{\mathbf{F}}^{lang}(v_1) - \hat{\mathbf{F}}^{lang}(v_2)\|_2, & v_1, v_2 \in \mathbf{M} \\ \mathcal{L}_{s3d} = \sum_{i=1}^N \sum_{k=1}^K \mathbf{f}_k^{sem} \log(\mathbf{f}_k^{sem} / \mathbf{f}_j^{sem}), & \end{cases} \quad (9)$$

where $\hat{\mathbf{F}}^{lang}$ is the rendered feature map and \mathbf{f}^{sem} is the feature attribute on the Gaussian. Such a strategy facilitates the language Gaussians are closely attached on the surface of corresponding objects.

Table 1. **2D Quantitative Results on LERF-OVS Dataset.** We report the open-vocabulary localization accuracy (%) and 2D semantic segmentation (IoU scores). LSeg [19] is a 2D open-vocabulary segmentation network, while other methods [20, 34] are per-scene optimized language field models. LSM [8] is the generalizable language Gaussian method. The **bold** denotes the best results.

Scene Type	LSeg [19]		LangSplat [34]		LangSurf [20]		LSM [8]		LangScene-X (Ours)	
	mAcc \uparrow	mIoU \uparrow	mAcc \uparrow	mIoU \uparrow	mAcc \uparrow	mIoU \uparrow	mAcc \uparrow	mIoU \uparrow	mAcc \uparrow	mIoU \uparrow
Teatime	65.22	30.58	30.43	15.81	35.57	18.82	44.46	19.62	78.91	45.07
Ramen	54.55	37.86	49.45	15.08	48.85	21.79	54.55	23.01	72.73	42.92
Kitchen	72.73	51.37	27.27	15.23	36.36	18.72	49.99	26.05	90.91	63.58
Overall	64.17	39.94	35.72	15.37	40.26	19.78	49.67	22.89	80.85	50.52

Table 2. **2D Quantitative Results on ScanNet Dataset.** We report the open-vocabulary localization accuracy (%) and 2D semantic segmentation (IoU scores). The **bold** denotes the best results.

Scene Type	LSeg [19]		LangSplat [34]		LangSurf [20]		LSM [8]		LangScene-X (Ours)	
	mAcc \uparrow	mIoU \uparrow	mAcc \uparrow	mIoU \uparrow	mAcc \uparrow	mIoU \uparrow	mAcc \uparrow	mIoU \uparrow	mAcc \uparrow	mIoU \uparrow
0085.00	63.64	28.71	42.73	28.98	68.18	30.47	67.65	39.09	95.45	51.68
0114.02	83.32	65.97	40.21	15.75	78.08	63.59	73.33	43.42	92.35	72.10
0616.00	84.62	66.96	46.15	11.35	68.34	44.30	78.65	56.41	97.85	70.71
0617.00	86.36	62.66	45.45	29.09	72.73	56.04	75.63	54.33	90.91	71.65
Overall	79.49	56.08	43.64	21.29	71.83	48.60	73.82	48.31	94.14	66.54

4. Experiment

4.1. Experiment Setup

Implementation Details. We conduct all experiments on 8 NVIDIA A800 (80G) GPUs. In the framework of LangScene-X, we choose CogVideoX [49] as the backbone architecture of our TriMap video diffusion, StableNormal [50] as our normal mapper \mathcal{D}_N , and SAM2 [37] as the semantic mapper \mathcal{D}_S . We first train TriMap video diffusion with the key-frame interpolation capability for one week on large-scale web data. Then we finetune it on 3D-consistent real data (*i.e.*, RealEstate-10K [55] and ACID [24]) with 2000 steps on the learning rate 1×10^{-5} . Next, we apply StableNormal to annotate 200 normal video clips of 3D scene data from RealEstate-10k and finetune TriMap video diffusion along with RGB videos with 800 steps. Finally, we apply SAM2 to annotate 300 clips of semantic video clips to fine-tune with RGB and normal videos with 1000 steps. The AdamW [29] optimizer is employed for optimization. All videos are center-cropped and resized to 720×480 resolution with 49 frames. Additionally, for the Language Quantized Compressor, we set the number of embeddings as $K = 2048$ and the channel as $D = 3$. Then, we train our model on large-scale open-world dataset COCO [23] with a batch size of 16 and 500,000 steps. We set loss weights $\lambda_1 = 1, \lambda_2 = 0.2, \lambda_3 = 0.5$ during the training. For Language Surface Fields training, we first train the Gaussian model with only RGB and normal loss

\mathcal{L}_{normal} for 5,000 steps, and then we utilize semantic losses (*e.g.*, $\mathcal{L}_{s2d}, \mathcal{L}_{s3d}$ for 5,000 steps).

Baseline and Metrics. To demonstrate our strong capability in building 3D language-embedded scenes from only sparse views, we compare our LangScene-X against four competitive baselines: LSeg [19], LangSplat [34], LangSurf [20], and LSM [8]. For fair comparison, we follow the dataset choice in LangSplat [34] and LangSurf [20] and conduct comparisons on LERF-OVS dataset [17] and ScanNet dataset [7]. The LERF dataset is an in-the-wild dataset captured by a handheld device, while ScanNet is a large scene dataset captured by RGB-D devices in complex indoor scenes. Each scene contains semantic labels at 3D level, making it suitable for 3D scene reconstruction and understanding tasks. For quantitative results, we report the standard metrics in semantic understanding, including open-vocabulary localization accuracy (mAcc) and semantic segmentation (mIoU scores).

4.2. Main Results

We conduct experiments on 2D open-vocabulary segmentation by querying the language features of each Gasussian with text on both LERF-OVS [16] and Scannet [7] datasets, as reported in Tab. 1 and Tab. 2. By comparing with existing state-of-the-art 3D language field techniques (*e.g.*, LangSplat, LangSurf), unified 3D representation method (*i.e.*, LSM), and open-vocabulary methods like LSeg, our method achieves superior performance in segmentation ac-

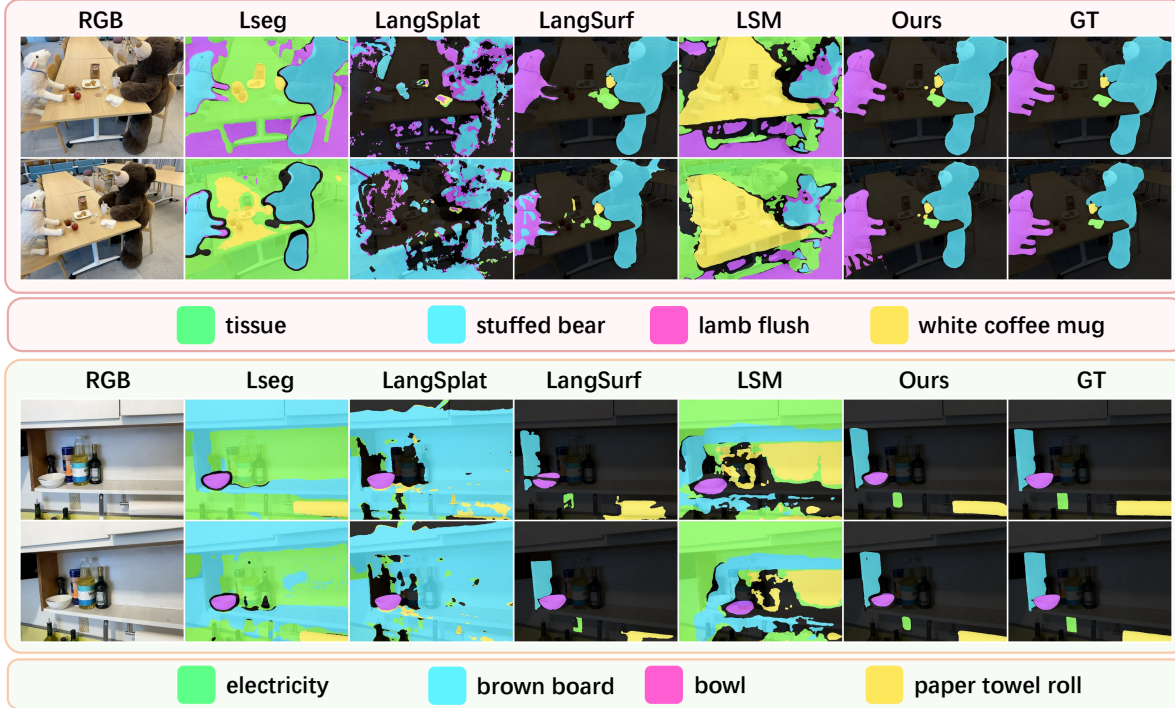


Figure 4. **2D Segmentation Results on LERF-OVS [17] Dataset.** Here, we showcase two cases (*i.e.*, Teatime, Kitchen) with multiple segmentation masks with text query. On the top, we display the rendered results of our method and other methods, along with the corresponding ground truth annotations. On the bottom, we present the images alongside the queried texts.

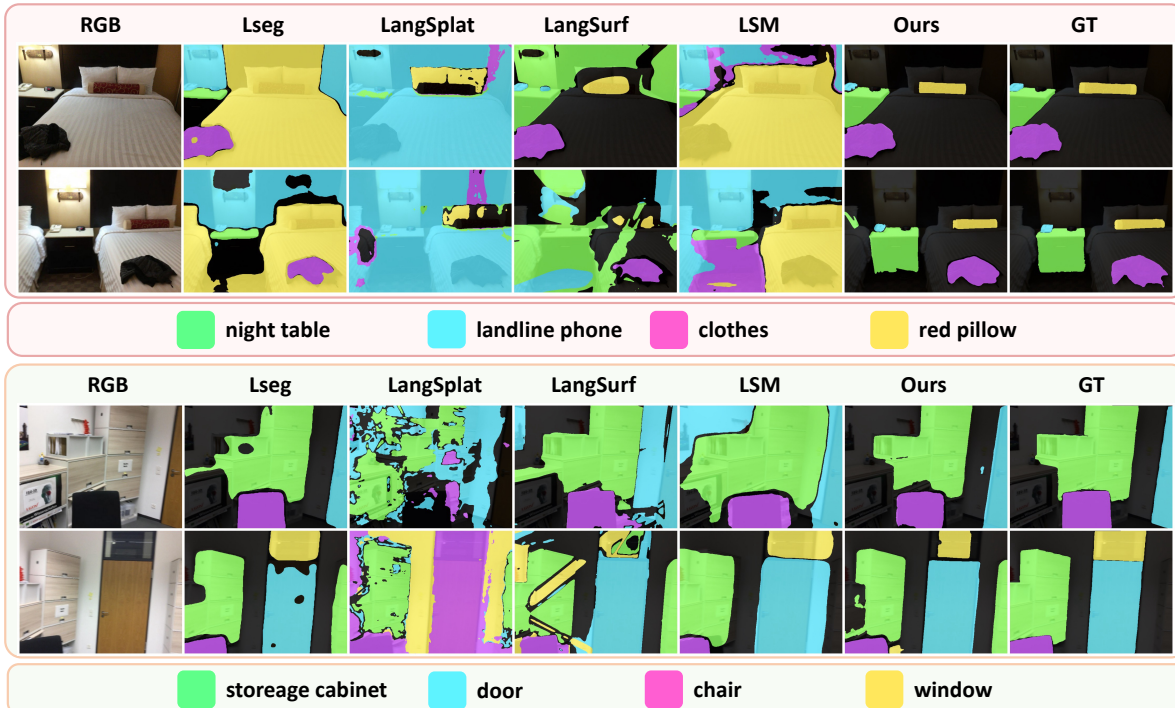


Figure 5. **2D Segmentation Results on Scannet [7] Dataset.** Here, we showcase two cases (*i.e.*, 0085.00, 0114.00) with multiple segmentation masks with text query. The masks predicted by ours contain more comprehensive regions and sharper boundaries than other methods, such as the “Cabinet” prompt, which also surpasses the GT masks.

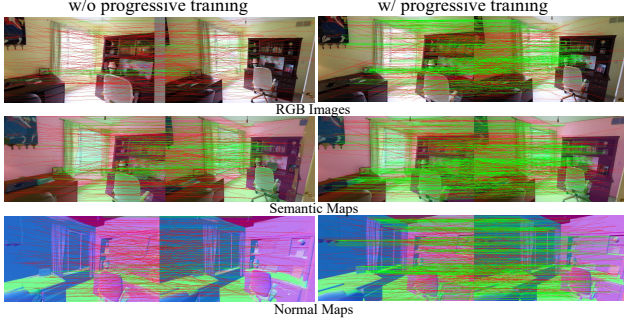


Figure 6. **Feature Matching** comparison between our method and vanilla video diffusion mdoel .

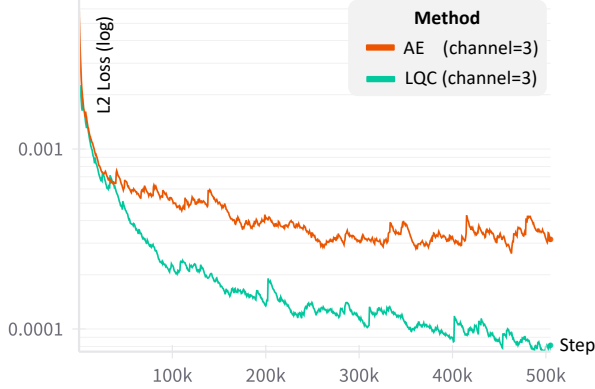


Figure 7. **Training Curve** comparison between our LQC and regular autoencoder technique.

curacy in both mIoU and mAcc metrics with a large margin, *i.e.*, a 10.58% in terms of mIoU and a 31.18% in terms of mAcc on LERF-OVS dataset. On the ScanNet dataset, the improvement upon the best existing method comes to 14.92% in terms of mIoU. The visualization of the 2D segmentation masks is shown in Fig. 4 and Fig. 5. As can be seen, our approach excels both optimized-based and feed-forward-based methods by segmenting objects with fine-grained boundaries. This demonstrates the capability of our method for embedding more generated multi-modality priors with 3D consistency into the 3D space.

4.3. Ablations

We conduct ablation experiments with our TriMap Video Diffusion and Language Quantized Compressor techniques. **TriMap Video Diffusion.** To further assess the geometric consistency of generated frames from TriMap video diffusion, we evaluate camera geometry alignment between frames. Specifically, we extract two frames at regular intervals from each video, creating pairs of two-view images. For each pair, we use a matching algorithm [32] to find corresponding points and use RANSAC [9] to filter out incorrect matches. Figure 6 shows the effectiveness of pro-

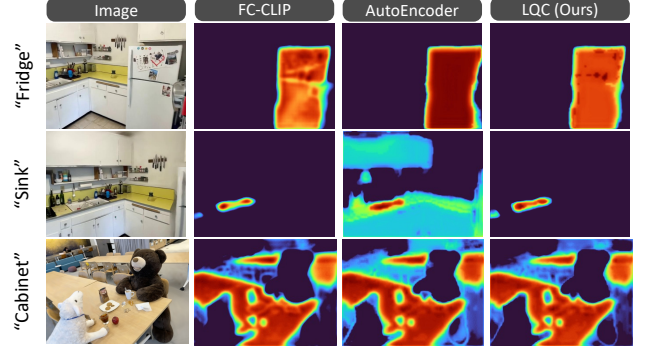


Figure 8. **Qualitative Comparison on LERF-OVS** [17]. We visualize text-query activation masks with various text prompts.

Table 3. **Ablations of proposed module and losses.** We perform ablations on both ScanNet (scene0085) [7] and LERF (Teaime) [17] with the segmentation metric mIoU.

Progressive Train	LQC	\mathcal{L}_{s2d}	\mathcal{L}_{s3d}	ScanNet	LERF
✗	✓	✓	✓	44.25	39.04
✓	✗	✓	✓	44.59	41.56
✓	✓	✗	✓	40.05	36.26
✓	✓	✓	✓	51.68	45.07

gressive training in TriMap video diffusion, which achieves more matched points.

Language Quantized Compressor. We visualize the training curve of our method and traditional autoencoder. As shown in Fig. 7, it is evident that our quantized technique demonstrates superior performance in both loss convergence rate and accuracy in terms of L2 loss (from $1e^{-3}$ to $1e^{-4}$). Moreover, we showcase the text-query activation masks in Fig. 8, where our method enable to perform sharper boundaries and more accurate activation scores within the query objects. This indicates the effectiveness of our discrete latent representation of LQC to preserve essential language properties during compression.

5. Conclusion

In this paper, we present LangScene-X, a generative framework that builds generalizable 3D language-embedded fields from only sparse views, which unify the information of reconstructing and understanding scenes in one video diffusion model. Specifically, we first train a TriMap video diffusion model through progressive knowledge integration, which can generate 3D consistent RGBs, normals, and semantic maps. Then we introduce a language quantized compressor to map high-dimensional language features into efficient feature representations. Finally, we reconstruct the language-embedded Gaussians by aligning the generated semantics onto the surface of 3D scenes. We believe LangScene-X provides a promising direction for 3D reconstruction and understanding through a generative paradigm.

References

- [1] Tomer Amit, Tal Shaharabany, Eliya Nachmani, and Lior Wolf. Segdiff: Image segmentation with diffusion probabilistic models. *arXiv preprint arXiv:2112.00390*, 2021. 2
- [2] Daichi Azuma, Taiki Miyanishi, Shuhei Kurita, and Motoaki Kawanabe. Scanqa: 3d question answering for spatial scene understanding. In *proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 19129–19139, 2022. 2
- [3] 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. 3
- [4] David Charatan, Sizhe Lester Li, Andrea Tagliasacchi, and Vincent Sitzmann. pixelsplat: 3d gaussian splats from image pairs for scalable generalizable 3d reconstruction. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 19457–19467, 2024. 3
- [5] Danpeng Chen, Hai Li, Weicai Ye, Yifan Wang, Weijian Xie, Shangjin Zhai, Nan Wang, Haomin Liu, Hujun Bao, and Guofeng Zhang. Pgsr: Planar-based gaussian splatting for efficient and high-fidelity surface reconstruction. *arXiv preprint arXiv:2406.06521*, 2024. 2
- [6] Yuedong Chen, Haofer Xu, Chuanxia Zheng, Bohan Zhuang, Marc Pollefeys, Andreas Geiger, Tat-Jen Cham, and Jianfei Cai. Mvsplat: Efficient 3d gaussian splatting from sparse multi-view images. In *European Conference on Computer Vision*, pages 370–386. Springer, 2024. 3
- [7] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5828–5839, 2017. 2, 6, 7, 8
- [8] Zhiwen Fan, Jian Zhang, Wenyan Cong, Peihao Wang, Renjie Li, Kairun Wen, Shijie Zhou, Achuta Kadambi, Zhangyang Wang, Danfei Xu, et al. Large spatial model: End-to-end unposed images to semantic 3d. *Advances in Neural Information Processing Systems*, 37:40212–40229, 2025. 3, 6
- [9] MA FISCHLER AND. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Commun. ACM*, 24(6):381–395, 1981. 8
- [10] Daniel Gordon, Aniruddha Kembhavi, Mohammad Rastegari, Joseph Redmon, Dieter Fox, and Ali Farhadi. Iqa: Visual question answering in interactive environments. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4089–4098, 2018. 2
- [11] Antoine Guédon and Vincent Lepetit. Sugar: Surface-aligned gaussian splatting for efficient 3d mesh reconstruction and high-quality mesh rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5354–5363, 2024. 2
- [12] 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. 3
- [13] Jun Hu, Zhang Chen, Zhong Li, Yi Xu, and Juyong Zhang. Sparselgs: Sparse view language embedded gaussian splatting. *arXiv preprint arXiv:2412.02245*, 2024. 2
- [14] Wenbo Hu, Xiangjun Gao, Xiaoyu Li, Sijie Zhao, Xiaodong Cun, Yong Zhang, Long Quan, and Ying Shan. Depthcrafter: Generating consistent long depth sequences for open-world videos. *arXiv preprint arXiv:2409.02095*, 2024. 3
- [15] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 2023. 2, 3
- [16] Justin Kerr, Chung Min Kim, Ken Goldberg, Angjoo Kanazawa, and Matthew Tancik. Lrf: Language embedded radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19729–19739, 2023. 6
- [17] Justin Kerr, Chung Min Kim, Ken Goldberg, Angjoo Kanazawa, and Matthew Tancik. Lrf: Language embedded radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19729–19739, 2023. 2, 3, 6, 7, 8
- [18] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. Segment anything. *arXiv:2304.02643*, 2023. 3
- [19] Boyi Li, Kilian Q Weinberger, Serge Belongie, Vladlen Koltun, and Rene Ranftl. Language-driven semantic segmentation. In *International Conference on Learning Representations*, 2022. 5, 6
- [20] Hao Li, Roy Qin, Zhengyu Zou, Diqi He, Bohan Li, Bingquan Dai, Dingwen Zhang, and Junwei Han. Langsurf: Language-embedded surface gaussians for 3d scene understanding. *arXiv preprint arXiv:2412.17635*, 2024. 2, 5, 6
- [21] Hao Li, Dingwen Zhang, Yalun Dai, Nian Liu, Lechao Cheng, Jingfeng Li, Jingdong Wang, and Junwei Han. Gpnerf: Generalized perception nerf for context-aware 3d scene understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21708–21718, 2024. 3
- [22] Jiahe Li, Jiawei Zhang, Xiao Bai, Jin Zheng, Xin Ning, Jun Zhou, and Lin Gu. Dngaussian: Optimizing sparse-view 3d gaussian radiance fields with global-local depth normalization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 20775–20785, 2024. 3
- [23] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer vision—ECCV 2014: 13th European conference, Zurich, Switzerland, September 6–12, 2014, proceedings, part v 13*, pages 740–755. Springer, 2014. 5, 6
- [24] Andrew Liu, Richard Tucker, Varun Jampani, Ameesh Makadia, Noah Snavely, and Angjoo Kanazawa. Infinite nature: Perpetual view generation of natural scenes from a single image. In *Proceedings of the IEEE/CVF International*

- Conference on Computer Vision*, pages 14458–14467, 2021. 6
- [25] Fangfu Liu, Chubin Zhang, Yu Zheng, and Yueqi Duan. Semantic ray: Learning a generalizable semantic field with cross-reprojection attention. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17386–17396, 2023. 2
- [26] Fangfu Liu, Wenqiang Sun, Hanyang Wang, Yikai Wang, Haowen Sun, Junliang Ye, Jun Zhang, and Yueqi Duan. Re-conx: Reconstruct any scene from sparse views with video diffusion model. *arXiv preprint arXiv:2408.16767*, 2024. 3, 4
- [27] Fangfu Liu, Hanyang Wang, Weiliang Chen, Haowen Sun, and Yueqi Duan. Make-your-3d: Fast and consistent subject-driven 3d content generation. In *European Conference on Computer Vision*, pages 389–406. Springer, 2024. 2
- [28] Fangfu Liu, Hanyang Wang, Shunyu Yao, Shengjun Zhang, Jie Zhou, and Yueqi Duan. Physics3d: Learning physical properties of 3d gaussians via video diffusion. *arXiv preprint arXiv:2406.04338*, 2024. 2
- [29] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 6
- [30] Hidenobu Matsuki, Riku Murai, Paul HJ Kelly, and Andrew J Davison. Gaussian splatting slam. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18039–18048, 2024. 2
- [31] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021. 2, 3
- [32] Pauline C Ng and Steven Henikoff. Sift: Predicting amino acid changes that affect protein function. *Nucleic acids research*, 31(13):3812–3814, 2003. 8
- [33] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4195–4205, 2023. 4
- [34] Minghan Qin, Wanhua Li, Jiawei Zhou, Haoqian Wang, and Hanspeter Pfister. Langsplat: 3d language gaussian splatting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20051–20060, 2024. 2, 3, 5, 6
- [35] Alec Radford, Jong Wook Kim, Chris Hallacy, A. Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *ICML*, 2021. 2, 3
- [36] Adam Rashid, Satvik Sharma, Chung Min Kim, Justin Kerr, Lawrence Yunliang Chen, Angjoo Kanazawa, and Ken Goldberg. Language embedded radiance fields for zero-shot task-oriented grasping. In *7th Annual Conference on Robot Learning*, 2023. 2
- [37] 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. 6
- [38] Rahul Ravishankar, Zeeshan Patel, Jathushan Rajasegaran, and Jitendra Malik. Scaling properties of diffusion models for perceptual tasks. *arXiv preprint arXiv:2411.08034*, 2024. 2, 4
- [39] Xinru Shan and Chaoning Zhang. Robustness of segment anything model (sam) for autonomous driving in adverse weather conditions. *arXiv preprint arXiv:2306.13290*, 2023. 2
- [40] Jiahao Shao, Yuanbo Yang, Hongyu Zhou, Youmin Zhang, Yujun Shen, Vitor Guizilini, Yue Wang, Matteo Poggi, and Yiyi Liao. Learning temporally consistent video depth from video diffusion priors. *arXiv preprint arXiv:2406.01493*, 2024. 4
- [41] Wenqiang Sun, Shuo Chen, Fangfu Liu, Zilong Chen, Yueqi Duan, Jun Zhang, and Yikai Wang. Dimensionx: Create any 3d and 4d scenes from a single image with controllable video diffusion. *arXiv preprint arXiv:2411.04928*, 2024. 3
- [42] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in neural information processing systems*, 30, 2017. 5
- [43] A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017. 3
- [44] 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. 4
- [45] Shuzhe Wang, Vincent Leroy, Yohann Cabon, Boris Chidlovskii, and Jerome Revaud. Dust3r: Geometric 3d vision made easy. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20697–20709, 2024. 3, 5
- [46] Yanmin Wu, Jiarui Meng, Haijie Li, Chenming Wu, Yahao Shi, Xinhua Cheng, Chen Zhao, Haocheng Feng, Errui Ding, Jingdong Wang, et al. Opengaussian: Towards point-level 3d gaussian-based open vocabulary understanding. *arXiv preprint arXiv:2406.02058*, 2024. 3
- [47] Jinbo Xing, Menghan Xia, Yong Zhang, Haoxin Chen, Wangbo Yu, Hanyuan Liu, Gongye Liu, Xintao Wang, Ying Shan, and Tien-Tsin Wong. Dynamicrafter: Animating open-domain images with video diffusion priors. In *European Conference on Computer Vision*, pages 399–417. Springer, 2024. 3
- [48] Jihan Yang, Shusheng Yang, Anjali W Gupta, Rilyn Han, Li Fei-Fei, and Saining Xie. Thinking in space: How multimodal large language models see, remember, and recall spaces. *arXiv preprint arXiv:2412.14171*, 2024. 2
- [49] 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. 3, 4, 6
- [50] Chongjie Ye, Lingteng Qiu, Xiaodong Gu, Qi Zuo, Yushuang Wu, Zilong Dong, Liefeng Bo, Yuliang Xiu, and Xiaoguang Han. Stablenormal: Reducing diffusion variance for stable and sharp normal. *ACM Transactions on Graphics (TOG)*, 43(6):1–18, 2024. 2, 6

- [51] Qihang Yu, Ju He, Xueqing Deng, Xiaohui Shen, and Liang-Chieh Chen. Convolutions die hard: Open-vocabulary segmentation with single frozen convolutional clip. *Advances in Neural Information Processing Systems*, 36:32215–32234, 2023. [5](#)
- [52] Wangbo Yu, Jinbo Xing, Li Yuan, Wenbo Hu, Xiaoyu Li, Zhipeng Huang, Xiangjun Gao, Tien-Tsin Wong, Ying Shan, and Yonghong Tian. Viewcrafter: Taming video diffusion models for high-fidelity novel view synthesis. *arXiv preprint arXiv:2409.02048*, 2024. [3](#), [4](#)
- [53] Junjie Zhang, Chenjia Bai, Haoran He, Wenke Xia, Zhigang Wang, Bin Zhao, Xiu Li, and Xuelong Li. Sam-e: leveraging visual foundation model with sequence imitation for embodied manipulation. *arXiv preprint arXiv:2405.19586*, 2024. [2](#)
- [54] Rui Zhao, Yuchao Gu, Jay Zhangjie Wu, David Junhao Zhang, Jia-Wei Liu, Weijia Wu, Jussi Keppo, and Mike Zheng Shou. Motiondirector: Motion customization of text-to-video diffusion models. In *European Conference on Computer Vision*, pages 273–290. Springer, 2024. [3](#)
- [55] 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. [6](#)