

MegaScenes: Scene-Level View Synthesis at Scale

Joseph Tung^{*1}, Gene Chou^{*1}, Ruojin Cai¹, Guandao Yang², Kai Zhang³,
Gordon Wetzstein², Bharath Hariharan¹, and Noah Snavely¹

¹ Cornell University
² Stanford University
³ Adobe Research

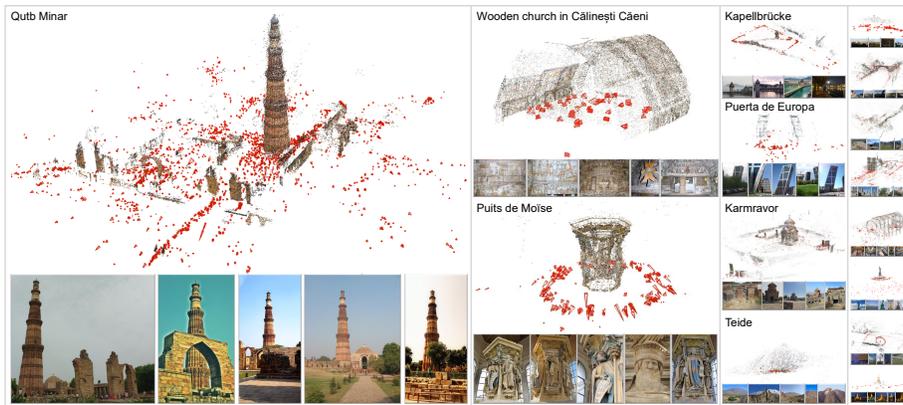


Fig. 1: The MegaScenes Dataset is an extensive collection of around 430k scenes, featuring over 100k structure-from-motion reconstructions and over 2 million registered images. MegaScenes includes a diverse array of scenes, such as minarets (e.g., Qutb Minar), building interiors (e.g., wooden church in Călinești Căeni), statues (e.g., Puits de Moïse), bridges (e.g., Kapellbrücke), towers (e.g., Puerta de Europa), religious buildings (e.g., Karmravor), and natural landscapes (e.g., Teide volcano). The images of these scenes are captured under varying conditions, including different times of day, various weather and illumination, and from different devices with distinct camera intrinsics.

Abstract. Scene-level novel view synthesis (NVS) is fundamental to many vision and graphics applications. Recently, pose-conditioned diffusion models have led to significant progress by extracting 3D information from 2D foundation models, but these methods are limited by the lack of scene-level training data. Common dataset choices either consist of isolated objects (Objaverse), or of object-centric scenes with limited pose distributions (DTU, CO3D). In this paper, we create a large-scale scene-level dataset from Internet photo collections, called MegaScenes, which contains over 100K structure from motion (SfM) reconstructions

* Equal contribution.

from around the world. Internet photos represent a scalable data source but come with challenges such as lighting and transient objects. We address these issues to further create a subset suitable for the task of NVS. Additionally, we analyze failure cases of state-of-the-art NVS methods and significantly improve generation consistency. Through extensive experiments, we validate the effectiveness of both our dataset and method on generating in-the-wild scenes. For details on the dataset and code, see our project page at <https://megascenes.github.io>.

Keywords: Novel view synthesis of scenes · Pose-conditioned diffusion models · Dataset of Internet photo collections

1 Introduction

Our vast visual experience enables us to look at a single view of a scene and infer what we cannot see. We can see a bridge from afar and imagine what it would be like to stand under it, or view the front of a church and guess what it looks like from other sides. Imagine a computer vision model that has similarly seen countless scenes: like humans, it can infer other views of a scene from a single image (i.e., it can perform *single-view novel-view synthesis*). Beyond connections with human vision, such a vision model would allow us to explore new AR/VR visualizations [67] or plan effectively in robotics [12, 68].

Current state-of-the-art methods on single-view novel-view synthesis (NVS) take 2D diffusion models trained on large internet datasets [39] and finetunes them on multiview images with camera poses. Concretely, these models map a reference image and a target pose to a target view [27, 56]. These methods successfully produced consistent novel views at an object level, as they were trained on object meshes. Unfortunately, attempts to generalize this approach to scenes [6, 41] by training on existing scene-level datasets [25, 37, 45, 67] were held back by the relatively small size and lack of diversity of these scene-level datasets. As such, current scene-level NVS techniques fail to match the consistency of object-level models and generalize to realistic, in-the-wild scenes.

To address the lack of diverse, scene-level data for training 3D-aware models, we create MegaScenes, a large-scale 3D dataset. MegaScenes builds on eight million free-to-use images sourced from Wikimedia Commons. We leverage structure from motion (SfM) to extract 3D structure from internet images at scale. In total, MegaScenes contains over 100K scene-level SfM reconstructions from around the world, along with associated data like captions, as well as the estimated relative poses of tens of millions of image pairs. Fig. 1 shows a few example scenes.

While we foresee a variety of 3D-related applications that could benefit from MegaScenes, such as pose estimation [50], feature matching [52], and reconstruction [53], in this paper we focus on NVS as a representative application. Following prior work in NVS [27, 41], our goal is to generate a plausible image at a target pose given only one reference image. Therefore, from MegaScenes we sample image pairs that have consistent lighting and visual overlap to create over 2 million training pairs. We validate MegaScenes’s effectiveness by finetuning

current state-of-the-art NVS models on our dataset, and find that new models perform significantly better on multiple dataset benchmarks.

In these experiments, we also identify and mitigate failure cases of existing methods by including additional conditioning that warps the input image to the target view [25]. While our method is simple and builds on existing approaches, it addresses fundamental issues in prior works, and we validate that it produces significantly more consistent and realistic results.

We show extensive experiments in Sec. 4 and a large collection of uncurated results in the supplement to demonstrate that our method and our training dataset yield NVS models that are effective across multiple benchmarks. We will release the dataset, code, and pretrained models.

2 Related Work

Datasets for 3D Learning. Datasets are the keystone of 3D learning. Recently, many 3D datasets have provided increasing amounts of data for tasks such as novel view synthesis, scene understanding, and 3D generation. Object-level 3D datasets like ShapeNet [8], CO3Dv2 [37], and DTU [45] have been extensively utilized in sparse view NVS [63] and 3D generation [69]. The emergence of larger-scale 3D object datasets like MVIImgNet [65] and Objaverse-XL [11] has enabled more generalizable models [16, 27, 51] for 3D reconstruction and generation. However, these datasets are confined to objects and do not extend to full scenes.

At the scene scale, existing datasets [7, 10, 19, 22, 25, 38, 61, 62, 67] have facilitated scene-level view synthesis and generation, but are often limited to a constrained set of categories, such as indoor scenes and drone shots of nature. DL3DV-10K [24] is concurrent work that aims to create a diverse and large-scale 3D scene dataset from videos, but features limited variation in camera poses.

In contrast, scene-level 3D datasets sourced from internet photos, such as MegaDepth [20], present a diverse distribution of camera poses and intrinsics, various lighting conditions and weather, different times of day, and transient objects and is widely applied in monocular depth estimation [2, 60] and learned feature matching [13, 23, 47]. However, MegaDepth is limited in scale to just 196 landmarks. Two more recent scene-level datasets include Google Landmarks v2 [57] and WikiScenes [59], which also gather images from Wikimedia Commons. However, Google Landmarks only focuses on 2D retrieval (no 3D information), and WikiScenes focuses on specific categories like cathedrals.

To address these limitations, MegaScenes incorporates diverse scene categories that include indoor, outdoor, natural scenes, and object-like scenes such as statues. It significantly extends the scale of 3D scene data, surpassing MegaDepth by several orders of magnitude, and includes 3D annotations of camera poses and reconstructions. Sourced from the Wikimedia Foundation, MegaScenes benefits from rich metadata and a wide distribution of illumination and camera poses. Our findings in novel view synthesis demonstrate that image diversity within the *same* scene enhances model generalization capabilities, highlighting MegaScenes’s value in advancing the field of 3D learning.

Novel View Synthesis from Sparse Views. Novel view synthesis (NVS) is the task of generating images from unseen views given some known images of a scene. When many input views are available, one can reconstruct an explicit 3D scene model, e.g., a neural radiance field [31] or 3D Gaussians [18]. However, given only sparse views (or just one), methods must rely on heuristic priors such as geometry smoothness [32, 48] or data priors [6, 46, 58, 63]. Recently, a popular line of work uses foundation generative models [39, 40] as prior knowledge. To work around the lack of 3D data and instead use 2D foundation models, [33, 49, 54, 58] generate 3D objects by enforcing rendered images from unseen viewpoints to agree with generative models. [5, 9, 64] explicitly extract multiview images by warping reference images and their depth given target poses following Liu *et al.* [25], and use an inpainting model to fill in missing regions. However, since these 2D generative models are not 3D-aware, methods can suffer from artifacts such as the multiface problem [33] or inconsistent geometries [9, 64].

A promising alternative is to leverage generative models that can perform novel-view synthesis conditioned on input view and change of camera poses [5, 6, 16, 17, 26–28, 34, 41, 44, 56, 69]. These methods produce consistent geometry given sparse or even single views without artifacts from 2D models, and can generalize to unseen scenarios thanks to their data priors. However, these works are generally trained on data with limited diversity, such as on object meshes [11] and object-centric scenes [37]. As we will show later, this limits their applicability to realistic, in-the-wild scenes. In this paper, we create both a dataset and a method that directly addresses scene-level novel view synthesis.

3 MegaScenes Dataset

In this section, we introduce the MegaScenes dataset, designed to capture a diverse range of geometries for large-scale scenes—plazas, buildings, interiors, and natural landmarks—using worldwide internet photos. We describe the dataset’s key characteristic features in Sec. 3.1. We detail the data collection and reconstruction pipeline of MegaScenes in Sec. 3.2 We provide dataset statistics in Sec. 3.3.

3.1 Dataset Characteristics

We describe several characteristics that highlight MegaScenes’s versatility for future vision tasks, including category, image, and 3D information.

Wikimedia Commons Categories as Scenes. Each scene in MegaScenes is based on single Wikimedia Commons category. Contributors from around the world have uploaded millions of images to Wikimedia Commons, and have organized images into representative groups. As shown in Fig. 3 we find that Wikimedia Commons categories depict scenes that are distributed across Earth, making it suitable as the foundation for a diverse dataset and future expandability.

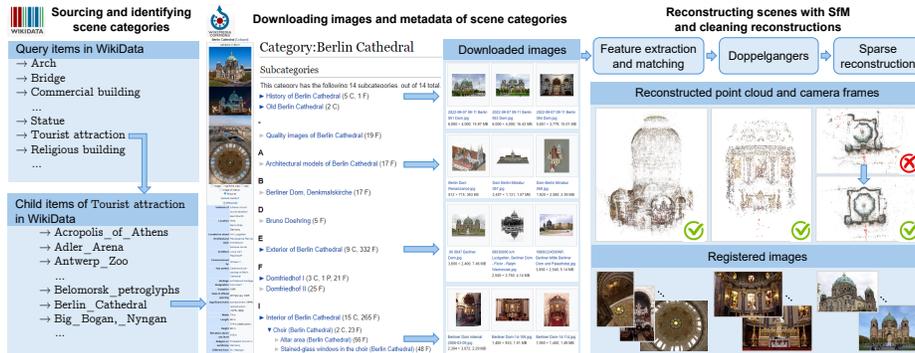


Fig. 2: MegaScenes curation pipeline. We first source and identify potential scene categories from WikiData. Subsequently, images and metadata for each scene category is downloaded. Finally, we reconstruct scenes using Structure from Motion (SfM) and clean them using the Doppelgangers [4] pipeline.

Images, Subcategorization, and Licensing. Images within a single scene are further classified into subcategories determined by Wikimedia Commons contributors. This enables future dataset applications to create subsets of data with greater granularity. This also proves to be helpful in cleaning the dataset, as described in the supplement on the dataset pipeline. Most importantly, like the similarly-sourced Google Landmarks v2 dataset [57], these images possess open content licenses or are in the public domain. Consequently, depending on the specific license, these images are free to reuse and alter for downstream tasks, so long as the original source is attributed.

3D Data. For each scene, we contribute SIFT [29] keypoints and descriptors, as well as calculated two-view geometries for pairs of images. We also contribute sparse point clouds and camera poses for a subset of scenes with ample image overlap. We use COLMAP [42] to compute this data.

Class Hierarchy. Similar to hierarchical extension of Google Landmarks v2 [36], the MegaScenes Dataset contains a hierarchy of class labels for each scene directly sourced from Wikidata. Wikidata is a large database of structured data connecting topics between Wikimedia Commons and Wikipedia. We use this class hierarchy to aid in dataset curation, as described in Sec. 3.2.

3.2 Dataset Curation

Fig. 2 depicts our dataset curation pipeline, which has three main steps: identifying scene categories, downloading images, and reconstructing scenes. We provide an overview of our pipeline below, and supply additional details in the supplement.

Our first goal is to identify Wikimedia Commons categories that may be considered as scenes. We take a top-down approach to identify scenes by utilizing the class hierarchy described in Sec. 3.1, as follows. First, we select several broad classes from Wikidata, such as “bridges” or “religious buildings”, that relate to collections of scenes. We choose these classes based on commonly seen places in everyday life. From these classes, we use the class hierarchy to identify Wikimedia Commons categories that are instances of these classes.

Next, we download all images associated with a Wikimedia Commons category that is identified as a scene, contingent on a subcategory filter we put in place to avoid downloading unrelated images. This filter is described in the supplement.

Lastly, for each scene, we run structure from motion on its corresponding collection of images using COLMAP [42] to produce sparse point clouds and camera poses. We use default parameters for feature extraction, vocabulary tree matching [43], and sparse reconstruction. We identify incorrect SfM reconstructions due to visual ambiguities (e.g., repeated patterns) by manual inspection guided by historically problematic scenes described in prior work [4, 14]. For these scenes, we run Doppelgangers [4] to get a corrected reconstruction.

3.3 Dataset Statistics

In total, MegaScenes consists of approximately 430K scenes derived from Wikimedia Commons. Across these categories, we download 9M images which results in over 30M image pairs with estimated two-view geometries. Around 80K of these scenes led to at least one sparse COLMAP reconstruction, resulting in over 100K reconstructions and 2M registered images. In these sparse reconstructions we triangulate 400 million 3D points, with a mean track length of 5 images and a mean of 8,700 observations per registered image. Similar to Google Landmarks [57], MegaScenes has a wide range of scenes, with as many as 18K images to as few as zero per scene. As shown by Fig. 3, our scenes covers a diverse set of classes ranging from buildings and outdoor spaces, to statues and streets.

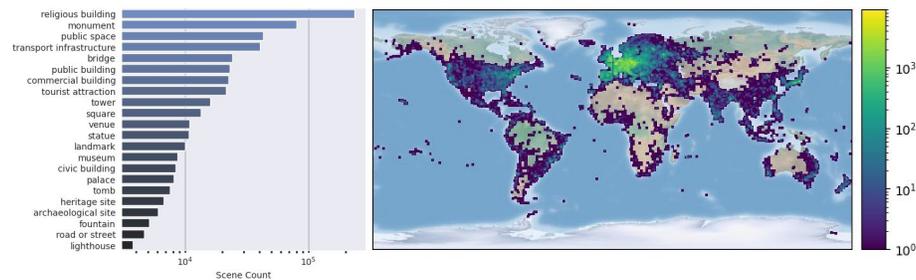


Fig. 3: Distribution of the MegaScenes Dataset. On the left, we depict the frequency of scenes grouped by Wikidata class. This includes only select classes with more than 3,500 scenes; note that a single scene may be an instance of multiple classes. On the right, we visualize the geospatial distribution of collected scenes worldwide.

4 MegaScenes Applied to Novel View Synthesis

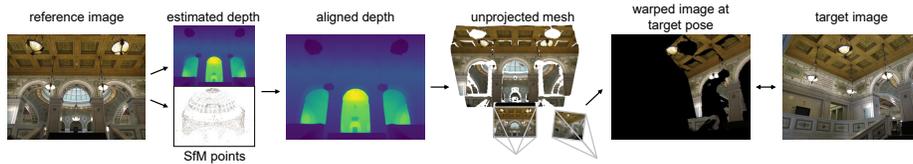


Fig. 4: We create over 2 million pairs of training images for novel view synthesis. Each pair contains relative pose and a warping from the reference to target image which we use for both training and evaluation. We align estimated monocular depths with sparse point clouds from COLMAP [42], and unproject the RGBD images to a mesh for viewpoint rendering. See Sec. 4.1 for details and Fig. 5 for more examples.

In this section, we explore MegaScenes on a representative application: novel view synthesis (NVS) from a single image. The goal is to take a reference image and generate a plausible image at a target pose that is consistent with the reference image. Following Sargent *et al.* [41], we train and evaluate on image pairs with pseudo-ground-truth relative poses obtained via SfM. In the supplement, we provide videos obtained through autoregressive generation.

We start with testing state-of-the-art novel-view synthesis models, namely Zero-1-to-3 [27] and ZeroNVS [41], on MegaScenes, and demonstrate that these approaches fail to generalize to in-the-wild scenes. We then improve these models in two ways. First, simply fine-tuning these methods on large numbers of training pairs from MegaScenes leads to dramatically improved results on both Internet photos of scenes and three out-of-domain datasets. Second, we observe that these fine-tuned models still demonstrate inconsistencies between the requested pose and the synthesized image. We show that by adding an additional conditioning image approximately warped from input view to target view, we improve pose consistency and novel view quality.

In Sec. 4.1, we describe our setup. Then, we show results of finetuning baseline models on MegaScenes in Sec. 4.2. In Sec. 4.3, we analyze failure cases of existing methods and propose our method to improve pose consistency. Finally, we evaluate our method on multiple datasets in Sec. 4.4 and Sec. 4.5.

4.1 Setup: Data Mining and Evaluation

Data Mining. We first identify a subset of image pairs from MegaScenes suitable for training novel view synthesis methods using two conditions. First, each pair should have similar lighting, since diffusion models operate on a pixel-wise loss. Using metadata, we find pairs of images taken within three hours of each other, as a proxy for lighting similarity. Second, we find pairs with sufficient visual overlap of at least 50 3D SfM points, so the model can learn view synthesis

based on visual cues. As shown in Fig. 5, we still observe both small and large pose changes with this threshold. Finally, we require that pairs have the same aspect ratio. Most previous works [39, 41] center crop images, but we find that many landmarks, such as statues, can have highly varied aspect ratios and lose information through center cropping. Thus, we resize the long side to 256 and pad the short side, to obtain images with size 256×256 .

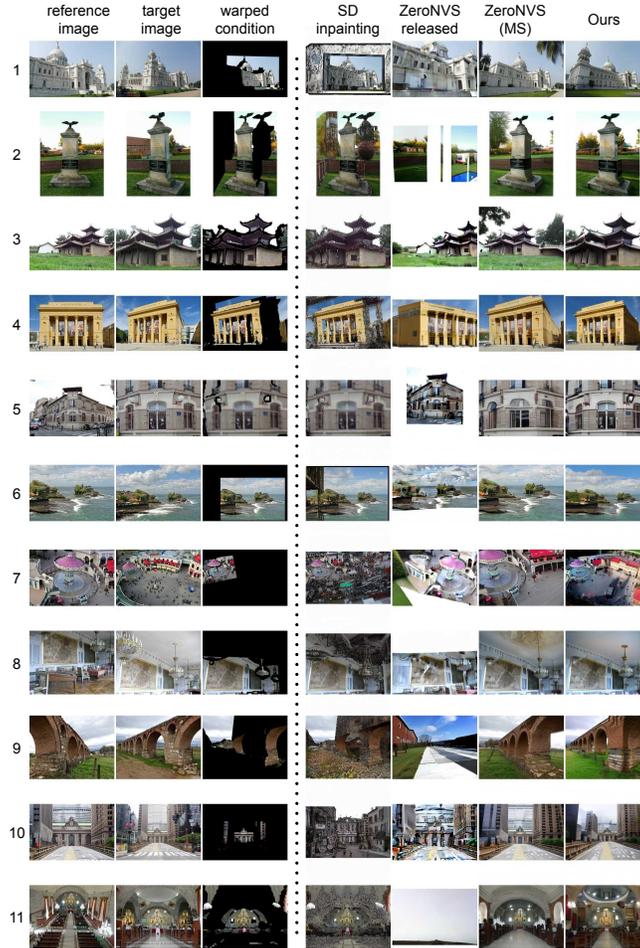


Fig. 5: We evaluate multiple baselines on MegaScenes, which contains diverse scenes, poses, and object compositions. Prior methods exhibit many failure modes in this challenging setting. Our method identifies and addresses these failure modes.

As a final check, we manually inspect all scenes and remove 298 scenes that we determine have too many occlusions in the majority of images; most of these

occlusions were people. In total, we obtain 2,086,036 pairs from 32,259 scenes and 475,277 unique images. We hold out 800 scenes that contain 51,240 pairs and 11,852 unique images. We form our validation set from the first 10,000 pairs, which we use to determine model convergence, and our test set with the remaining 41,240 pairs, which we use to report numbers.

Evaluation metrics. We evaluate each method using standard reconstruction and generation metrics. For reconstruction, we calculate LPIPS [66], PSNR, and SSIM [55]. LPIPS measures perceptual similarity, while PSNR and SSIM operate mainly on a pixel basis. However, generative models should only be expected to remain consistent in the target image where pixels from the reference image are present, and remain free to generate diverse samples, which could mean a lower reconstruction score. Thus, we propose “masked” versions of these metrics. We warp the input view to the target view using the target relative pose. Only pixels in the reference are present in the warped image; we only compare the copied pixels to the same location of the generated images. Our pipeline to create warpings is shown in Fig. 4. First, we use Depth-Anything [60] to estimate the reference image’s monocular depth. We project the COLMAP sparse point cloud to the reference image’s coordinate frame to obtain the ground-truth sparse depth, and use RANSAC to align the two. With this aligned dense depth, we unproject the reference RGBD to a mesh and render it from the target pose. For generative metrics, we use FID [15] and KID [3]. Both assess generated image quality by comparing their feature distributions to those of real images. Lower scores indicate that generated images are more similar to real images.

In general, we find LPIPS, FID, and KID reliable metrics for assessing *quality* and *realism* in generated images. We find Masked PSNR and Masked SSIM reliable for assessing *consistency*, i.e. whether generated images follow the target pose and retain details from the reference images. Still, we encourage readers to view qualitative results for a more comprehensive understanding.

4.2 Finetuning Pose-Conditioned Models on MegaScenes

Zero-1-to-3 [27] is finetuned from Stable Diffusion on Objaverse [11]. ZeroNVS [41] is finetuned from Zero-1-to-3 on CO3D [37], ACID [25], and RealEstate10K [67]. Our goal is to evaluate whether finetuning these models on MegaScenes improves generalization to in-the-wild scenes.

Finetuning details. Zero-1-to-3 conditions on poses in spherical coordinates, which is only suitable for objects placed in a canonical coordinate frame. Thus, we condition poses based on ZeroNVS, which flattens the extrinsic matrix and field of view as input to cross-attention, for both models. The scale of translation is determined by the 20th quantile of the depth of the reference image [41]. Additionally, both models concatenate the target and reference images and provide the reference image’s CLIP [35] embedding to cross-attention so that the output remains consistent with the reference. We compare released and finetuned

models to verify whether our dataset improves generalization to in-the-wild scenes. We provide training details in the supplement.

Table 1: We evaluate whether models trained on MegaScenes generalize to in-the-wild scenes. *Zero-1-to-3* / *ZeroNVS (released)* are released checkpoints. We finetune both on MegaScenes (models denoted with *(MS)*). *SD-inpainting* uses image warping and a pretrained diffusion inpainting model without finetuning, following the setup in [9, 64]. Our method takes warped images as input, and we condition with and without (*w/o ext*) the extrinsic matrix. \uparrow means higher is better and \downarrow means lower is better.

	LPIPS(\downarrow)	PSNR(\uparrow)	SSIM(\uparrow)	Masked LPIPS (\downarrow)	Masked PSNR (\uparrow)	Masked SSIM (\uparrow)	FID(\downarrow)	KID(\downarrow)
<i>Pose-Conditioned (Sec. 4.2)</i>								
Zero-1-to-3 (released)	0.5476	9.0896	0.2413	0.2777	14.132	0.6320	86.892	0.0634
ZeroNVS (released)	0.6156	7.4711	0.1508	0.3229	11.041	0.5421	69.097	0.0487
Zero-1-to-3 (MS)	0.4289	12.159	0.3665	0.1811	19.952	0.7286	9.7835	0.0023
ZeroNVS (MS)	0.3857	12.900	0.4005	0.1572	20.713	0.7534	9.8382	0.0024
<i>Warp-Conditioned (Sec. 4.3)</i>								
SD-inpainting	0.4245	12.358	0.3923	0.1283	24.377	0.8005	38.484	0.0242
Ours w/o ext	0.3534	13.310	0.4328	0.1297	22.609	0.7819	12.010	0.0041
<i>Warp + Pose (Sec. 4.3)</i>								
Ours	0.3444	13.397	0.4446	0.1256	22.483	0.7842	11.580	0.0040

Results. We show qualitative results in Fig. 5 and quantitative results in Tab. 1. We denote checkpoints released by authors with *(released)* and models finetuned on MegaScenes with *(MS)*. Additional results are in the supplement.

Zero-1-to-3 (released) and *ZeroNVS (released)* are both unable to generalize to internet photos. They produce unrealistic images with incorrect poses. We note that the former outperforms the latter in numbers, but upon inspecting qualitative results (see supplement) we observe that *Zero-1-to-3 (released)* tends to return the reference image. Finetuning on MegaScenes significantly improves results of both models, seen in the metrics of the *(MS)* models and the qualitative comparisons between *ZeroNVS (released)* and *ZeroNVS (MS)*.

We also validate that MegaScenes is suitable for the task of scene-level NVS. *Zero-1-to-3 (MS)* outperforms *ZeroNVS (released)* even though both are finetuned from *Zero-1-to-3*'s released checkpoint; one is trained on MegaScenes and the other on CO3D [37], ACID [25], and RealEstate10K [67].

ZeroNVS (MS) shows the best performance among these four models. From Fig. 5, we see that it produces realistic images, and the generated images clearly attempt to follow the desired pose. However, many images produced by *ZeroNVS (MS)* are still inaccurate. The positions of the islands (row 6), bridge (row 9), and building (row 10) are slightly different than in the target image, and when there is larger zoom, such as in rows 3, 4, and 7, the model fails to interpret the scale properly. Next, we address these issues.

4.3 Improving Pose Consistency with Warp Conditioning

ZeroNVS [41] conditions their model on the flattened extrinsic matrix, which is not a very intuitive pose condition; the model must learn a spatial transformation without visual cues. Furthermore, the translation scale is ambiguous since scenes cannot be canonicalized to a fixed coordinate frame. The authors in the original paper run a grid search on each scene to manually determine a scene scale during inference.

Since our goal is to generalize to in-the-wild scenes, we automatically determine scene scale using on monocular depth estimation when testing ZeroNVS. This leads to inaccurate poses especially with large zooming effects.

Our insight is that the warped image (Fig. 4) encodes pose by visualizing how pixels are supposed to move, and is directly aligned with the scene scale. On our training and evaluation datasets, the scale is based on 3D SfM points. Given a single image, we can determine the scene scale from estimated monocular depth and use the same extrinsics for conditioning and warping the image for a consistent scale. We show this setting in the supplement by generating videos given a single image. Thus, we concatenate the warped image with the input target and reference images, and observe significant improvements in pose accuracy.

However, using only the warped image as pose condition leads to two problems. First, inaccurate depth, which can be common in noisy, in-the-wild scenes, can cause the model to fail. Furthermore, with guidance only through the movement of 2D pixels, the model seems unable to interpret 3D structure at times.

Therefore, we also condition on the extrinsic matrix following ZeroNVS. We show qualitative results of these design choices in Fig. 6. In rows 1 and 4, *Ours* demonstrate better 3D consistency compared to *Ours w/o ext*, which is the model trained without conditioning on the extrinsic matrix, including creating the separating wall

and a complete building, respectively. In rows 2 and 3, there is little information in the warped condition due to inaccurate depth, and the generated results contain either unwanted objects or objects at inaccurate locations.

Using the warped image as a condition for view synthesis was first proposed by Liu *et al.* [25], then adapted by a recent line of work that uses a Stable Diffusion inpainting model without any finetuning [9, 64]. The premise is that a foundation model can generalize to any domain without forgetting. We therefore also compare to this baseline, denoted *SD-inpainting*. Using the preprocessed warpings, we set empty pixels as the mask to inpaint. This method, however, demonstrates lack

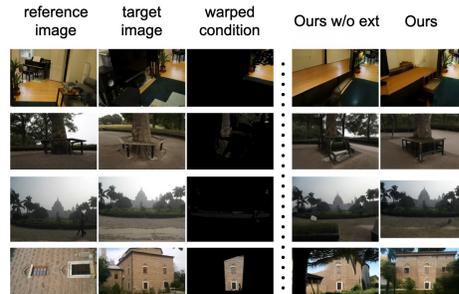


Fig. 6: We compare results with and without conditioning on the extrinsic matrix. The extrinsic matrix ensures valid outputs and produces more consistent 3D geometry.

of 3D understanding. For instance, the inpainting model frequently interprets a scene as a picture frame. Furthermore, the inpainting model is only trained on large, uniform masks, and produces artifacts when operating on fine-grained masks, which come with arbitrary poses. In contrast, the masks can be of arbitrary size, and our diffusion model has the freedom to remove occlusions and create plausible images. These results suggest that finetuning on 3D data is essential for zero-shot novel view synthesis.

In the following subsections, we evaluate all methods on MegaScenes as well as three out-of-domain datasets. While our method is simple and builds on existing methods, it addresses the fundamental issues of prior works, and we validate that it produces significantly more consistent and realistic results. We show a large collection of uncurated results in the supplement and demonstrate that our method is effective across a variety of diverse scenes.

4.4 Evaluation on MegaScenes

We first evaluate on MegaScenes’ test set, which consists of in-the-wild scenes from Internet photos. We show quantitative results in Tab. 1 and qualitative results in Fig. 5. Our method produces images closest to the desired pose, while being realistic and visually consistent with the reference. Compared to *ZeroNVS (MS)*, our method places objects according to cues in the warped image, which not only leads to more accurate positioning (and better reconstruction metrics), but also more detail (structure in row 5, statues at the end of the hall in row 11). Although our method also has higher FID and KID scores than the pose-conditioned models because the warped images add constraints to the generation process, we do not visually observe a degradation in image quality.

SD-inpainting has the best reconstruction metrics, as it faithfully returns the pixels in the warped condition. However, it does not understand 3D geometry, evident in the inconsistent generations, such as interpreting a scene as a picture frame (row 1). The inpainting artifacts also lead to unrealistic images, resulting in high LPIPS, FID, and KID scores. Our method avoids these issues while taking advantage of strong position cues of the warped condition.

4.5 Evaluation on Datasets from Different Domains

Next, we evaluate on DTU [45], Mip-NeRF 360 [1], and RealEstate10K (Re10K) [67], datasets commonly used for evaluating novel view synthesis. We expect MegaScenes to be sufficiently diverse such that models trained on it can generalize even to specific domains. We obtain image pairs and warpings from all three datasets. In total, we obtain 2,850 evaluation pairs from DTU, 15,682 pairs from Mip-NeRF 360, and 644 pairs from Re10K. We describe our data setup in the supplement.

We show results in Tab. 2 and Fig. 7. We see a similar trend as in Sec. 4.4. All metrics and qualitative results improve when trained on MegaScenes. *ZeroNVS (released)* was trained on object-centric datasets similar to DTU and Mip-NeRF 360, and directly trained on Re10K, but *Zero-1-to-3 (MS)* significantly outperforms it on all datasets. This validates that MegaScenes covers a wide variety of domains.

Table 2: We evaluate whether models trained on MegaScenes generalize to other data domains. The models and metrics are the same as in Tab. 1. We test on DTU, Mip-NeRF 360, and RealEstate10K.

DTU	LPIPS(↓)	PSNR(↑)	SSIM(↑)	Masked LPIPS (↓)	Masked PSNR (↑)	Masked SSIM (↑)	FID(↓)	KID(↓)
<i>Pose-Conditioned (Sec. 4.2)</i>								
Zero-1-to-3 (released)	0.5647	6.8720	0.2100	0.2592	12.628	0.6609	128.93	0.0297
ZeroNVS (released)	0.6476	5.7992	0.1113	0.3193	9.7005	0.5517	159.96	0.0352
Zero-1-to-3 (MS)	0.5158	7.6367	0.2755	0.2080	13.311	0.7014	101.94	0.0223
ZeroNVS (MS)	0.4833	8.0191	0.3066	0.1908	13.515	0.7152	87.406	0.0158
<i>Warp-Conditioned (Sec. 4.3)</i>								
SD-inpainting	0.4951	9.9463	0.3688	0.1283	22.656	0.8333	214.42	0.1067
Ours w/o ext	0.4113	8.8473	0.3878	0.1385	16.631	0.7924	92.284	0.0193
<i>Warp + Pose (Sec. 4.3)</i>								
Ours	0.3995	8.7953	0.3930	0.1357	16.593	0.7916	85.959	0.0163
Mip-NeRF 360	LPIPS(↓)	PSNR(↑)	SSIM(↑)	Masked LPIPS (↓)	Masked PSNR (↑)	Masked SSIM (↑)	FID(↓)	KID(↓)
<i>Pose-Conditioned</i>								
Zero-1-to-3 (released)	0.5258	10.720	0.2865	0.1621	16.299	0.8864	171.21	0.1126
ZeroNVS (released)	0.6685	6.9993	0.1240	0.2312	10.890	0.7670	137.04	0.0537
Zero-1-to-3 (MS)	0.4429	12.921	0.3828	0.0307	29.441	0.9697	67.645	0.0163
ZeroNVS (MS)	0.4057	13.780	0.4122	0.1369	24.909	0.8219	60.677	0.0139
<i>Warp-Conditioned</i>								
SD-inpainting	0.4557	12.922	0.3996	0.1212	27.455	0.8488	150.11	0.0792
Ours w/o ext	0.3944	13.667	0.4279	0.1237	25.884	0.8344	70.684	0.0193
<i>Warp + Pose</i>								
Ours	0.3807	14.056	0.4406	0.1150	26.196	0.8422	64.406	0.0142
RE10K	LPIPS(↓)	PSNR(↑)	SSIM(↑)	Masked LPIPS (↓)	Masked PSNR (↑)	Masked SSIM (↑)	FID(↓)	KID(↓)
<i>Pose-Conditioned</i>								
Zero-1-to-3 (released)	0.4050	11.632	0.4384	0.2732	14.079	0.6400	160.20	0.0725
ZeroNVS (released)	0.4563	9.4869	0.3527	0.3078	11.456	0.5565	123.01	0.0352
Zero-1-to-3 (MS)	0.2722	14.638	0.5697	0.1510	21.241	0.7637	68.908	0.0024
ZeroNVS (MS)	0.2053	16.015	0.6304	0.1176	20.609	0.8070	61.117	0.0024
<i>Warp-Conditioned</i>								
SD-inpainting	0.2694	15.541	0.6429	0.0929	29.056	0.8719	118.94	0.0396
Ours w/o ext	0.1922	16.105	0.6267	0.1109	23.147	0.7985	66.770	0.0057
<i>Warp + Pose</i>								
Ours	0.1774	17.224	0.6661	0.0942	24.259	0.8315	60.013	0.0023

Again, our method produces images closest to the desired pose. *ZeroNVS (MS)* mostly follows the pose condition, but the positions of the objects are less accurate. This is obvious in DTU where we can visually match the corners of the objects. Direct visual cues in the warped condition allow our model to preserve structure even in challenging cases, such as the bicycle in row 6.

As a side note, we would like to point out that *ZeroNVS (released)* appears to perform worse than shown in their original paper because of different testing settings. The original paper evaluates results only after SDS [33] optimization, which filters out noise and samples the mode of the diffusion outputs. Additionally, the authors run a grid search to manually determine scene scale. We were able to reproduce results shown in the original paper, but here we present feed-forward results without optimization or manual tuning for fair comparison.

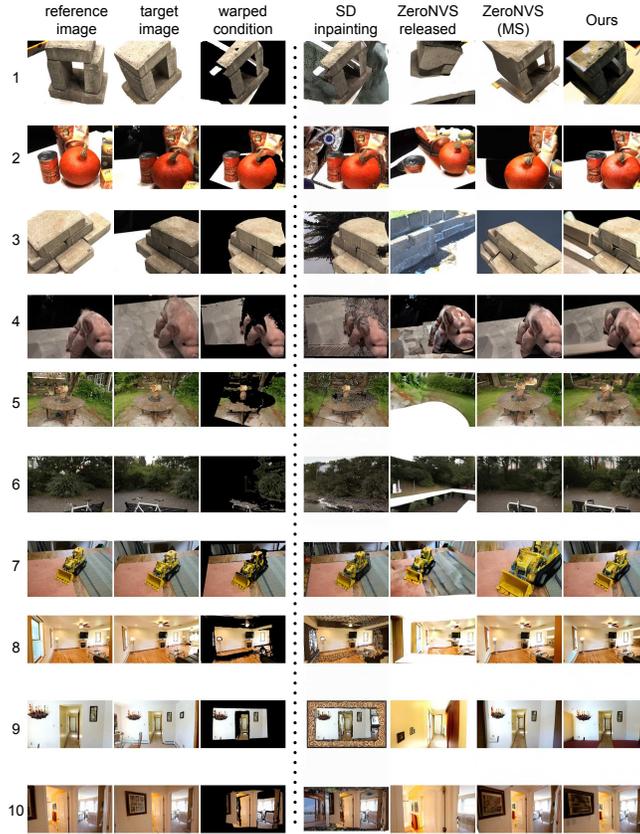


Fig. 7: We evaluate multiple models on DTU (rows 1-4), Mip-NeRF 360 (5-7), and Re10K (8-10). Models trained on MegaScenes are able to generalize to these datasets.

5 Conclusion

We present MegaScenes, a general large-scale 3D dataset, and analyze its impact on scene-level novel view synthesis. We find that finetuning NVS methods on MegaScenes significantly improves synthesis quality, which validates the dataset’s uses. We also improve existing methods and observe increased pose accuracy.

Regarding limitations and future work, on the task of NVS, we use a fraction of our data (475K out of 2M images) and a subset of data types (we did not use text captions). We would like to expand MegaScenes to applications that leverage the full dataset. Our NVS method also comes with limitations. It relies on warped images for conditioning and is impacted by erroneous depth estimation. Also, it cannot handle large camera motions such as behind a scene. Finally, we bypass lighting by sampling based on metadata, but we could incorporate lighting conditions [21, 30] in the future.

Acknowledgments

We thank Brandon Li for building the COLMAP web viewer. This work was funded in part by the National Science Foundation (IIS-2008313, IIS-2211259, IIS-2212084). Gene Chou was funded by an NSF Graduate Research Fellowship.

References

1. Barron, J.T., Mildenhall, B., Verbin, D., Srinivasan, P.P., Hedman, P.: Mip-nerf 360: Unbounded anti-aliased neural radiance fields. CVPR (2022)
2. Bhat, S.F., Birkel, R., Wofk, D., Wonka, P., Müller, M.: Zoedepth: Zero-shot transfer by combining relative and metric depth (2023). <https://doi.org/10.48550/ARXIV.2302.12288>, <https://arxiv.org/abs/2302.12288>
3. Bińkowski, M., Sutherland, D.J., Arbel, M., Gretton, A.: Demystifying mmd gans. arXiv preprint arXiv:1801.01401 (2018)
4. Cai, R., Tung, J., Wang, Q., Averbuch-Elor, H., Hariharan, B., Snavely, N.: Doppelgangers: Learning to disambiguate images of similar structures. In: ICCV (2023)
5. Cai, S., Chan, E.R., Peng, S., Shahbazi, M., Obukhov, A., Van Gool, L., Wetzstein, G.: Diffdreamer: Towards consistent unsupervised single-view scene extrapolation with conditional diffusion models. In: ICCV (2023)
6. Chan, E.R., Nagano, K., Chan, M.A., Bergman, A.W., Park, J.J., Levy, A., Aittala, M., De Mello, S., Karras, T., Wetzstein, G.: Generative novel view synthesis with 3d-aware diffusion models. arXiv preprint arXiv:2304.02602 (2023)
7. Chang, A., Dai, A., Funkhouser, T., Halber, M., Niessner, M., Savva, M., Song, S., Zeng, A., Zhang, Y.: Matterport3d: Learning from rgb-d data in indoor environments. arXiv preprint arXiv:1709.06158 (2017)
8. Chang, A.X., Funkhouser, T., Guibas, L., Hanrahan, P., Huang, Q., Li, Z., Savarese, S., Savva, M., Song, S., Su, H., et al.: Shapenet: An information-rich 3d model repository. arXiv preprint arXiv:1512.03012 (2015)
9. Chung, J., Lee, S., Nam, H., Lee, J., Lee, K.M.: Luciddreamer: Domain-free generation of 3d gaussian splatting scenes. arXiv preprint arXiv:2311.13384 (2023)
10. Dai, A., Chang, A.X., Savva, M., Halber, M., Funkhouser, T., Niessner, M.: Scannet: Richly-annotated 3d reconstructions of indoor scenes. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 5828–5839 (2017)
11. Deitke, M., Liu, R., Wallingford, M., Ngo, H., Michel, O., Kusupati, A., Fan, A., Laforte, C., Voleti, V., Gadre, S.Y., VanderBilt, E., Kembhavi, A., Vondrick, C., Gkioxari, G., Ehsani, K., Schmidt, L., Farhadi, A.: Objaverse-xl: A universe of 10m+ 3d objects. arXiv preprint arXiv:2307.05663 (2023)
12. Du, Y., Zhang, Y., Yu, H.X., Tenenbaum, J.B., Wu, J.: Neural radiance flow for 4d view synthesis and video processing. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (2021)
13. Edstedt, J., Bökman, G., Wadenbäck, M., Felsberg, M.: DeDoDe: Detect, Don't Describe — Describe, Don't Detect for Local Feature Matching. In: 2024 International Conference on 3D Vision (3DV). IEEE (2024)
14. Heinly, J., Dunn, E., Frahm, J.M.: Recovering Correct Reconstructions from Indistinguishable Geometry. In: International Conference on 3D Vision (3DV) (2014)
15. Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S.: Gans trained by a two time-scale update rule converge to a local nash equilibrium. Advances in neural information processing systems **30** (2017)

16. Hong, Y., Zhang, K., Gu, J., Bi, S., Zhou, Y., Liu, D., Liu, F., Sunkavalli, K., Bui, T., Tan, H.: Lrm: Large reconstruction model for single image to 3d. arXiv preprint arXiv:2311.04400 (2023)
17. Kant, Y., Siarohin, A., Vasilkovsky, M., Guler, R.A., Ren, J., Tulyakov, S., Gilitschenski, I.: invs : Repurposing diffusion inpainters for novel view synthesis. In: SIGGRAPH Asia 2023 Conference Papers (2023)
18. Kerbl, B., Kopanas, G., Leimkühler, T., Drettakis, G.: 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics* **42**(4) (2023)
19. Li, Y., Jiang, L., Xu, L., Xiangli, Y., Wang, Z., Lin, D., Dai, B.: Matrixcity: A large-scale city dataset for city-scale neural rendering and beyond. arXiv e-prints pp. arXiv-2308 (2023)
20. Li, Z., Snavely, N.: Megadepth: Learning single-view depth prediction from internet photos. In: *Computer Vision and Pattern Recognition (CVPR)* (2018)
21. Li, Z., Xian, W., Davis, A., Snavely, N.: Crowdsampling the plenoptic function. In: *European Conference on Computer Vision*. pp. 178–196. Springer (2020)
22. Li, Z., Yu, T.W., Sang, S., Wang, S., Song, M., Liu, Y., Yeh, Y.Y., Zhu, R., Gundavarapu, N., Shi, J., et al.: Openrooms: An open framework for photorealistic indoor scene datasets. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. pp. 7190–7199 (2021)
23. Lindenberger, P., Sarlin, P.E., Pollefeys, M.: LightGlue: Local Feature Matching at Light Speed. In: *ICCV* (2023)
24. Ling, L., Sheng, Y., Tu, Z., Zhao, W., Xin, C., Wan, K., Yu, L., Guo, Q., Yu, Z., Lu, Y., et al.: D13dv-10k: A large-scale scene dataset for deep learning-based 3d vision. arXiv preprint arXiv:2312.16256 (2023)
25. Liu, A., Tucker, R., Jampani, V., Makadia, A., Snavely, N., Kanazawa, A.: Infinite nature: Perpetual view generation of natural scenes from a single image. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)* (October 2021)
26. Liu, M., Xu, C., Jin, H., Chen, L., Varma T, M., Xu, Z., Su, H.: One-2-3-45: Any single image to 3d mesh in 45 seconds without per-shape optimization. *Advances in Neural Information Processing Systems* **36** (2024)
27. Liu, R., Wu, R., Van Hoorick, B., Tokmakov, P., Zakharov, S., Vondrick, C.: Zero-1-to-3: Zero-shot one image to 3d object. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 9298–9309 (2023)
28. Liu, Y., Lin, C., Zeng, Z., Long, X., Liu, L., Komura, T., Wang, W.: Syncdreamer: Generating multiview-consistent images from a single-view image. arXiv preprint arXiv:2309.03453 (2023)
29. Lowe, D.G.: Distinctive image features from scale-invariant keypoints. *IJCV* **60**(2), 91–110 (2004)
30. Martin-Brualla, R., Radwan, N., Sajjadi, M.S.M., Barron, J.T., Dosovitskiy, A., Duckworth, D.: NeRF in the Wild: Neural Radiance Fields for Unconstrained Photo Collections. In: *CVPR* (2021)
31. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM* **65**(1), 99–106 (2021)
32. Niemeyer, M., Barron, J.T., Mildenhall, B., Sajjadi, M.S., Geiger, A., Radwan, N.: Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 5480–5490 (2022)
33. Poole, B., Jain, A., Barron, J.T., Mildenhall, B.: Dreamfusion: Text-to-3d using 2d diffusion. arXiv preprint arXiv:2209.14988 (2022)

34. Qian, G., Mai, J., Hamdi, A., Ren, J., Siarohin, A., Li, B., Lee, H.Y., Skorokhodov, I., Wonka, P., Tulyakov, S., et al.: Magic123: One image to high-quality 3d object generation using both 2d and 3d diffusion priors. arXiv preprint arXiv:2306.17843 (2023)
35. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: International conference on machine learning. pp. 8748–8763. PMLR (2021)
36. Ramzi, E., Audebert, N., Rambour, C., Araujo, A., Bitot, X., Thome, N.: Optimization of Rank Losses for Image Retrieval. In: In submission to: IEEE Transactions on Pattern Analysis and Machine Intelligence (2023)
37. Reizenstein, J., Shapovalov, R., Henzler, P., Sbordone, L., Labatut, P., Novotny, D.: Common objects in 3d: Large-scale learning and evaluation of real-life 3d category reconstruction. In: International Conference on Computer Vision (2021)
38. Roberts, M., Ramapuram, J., Ranjan, A., Kumar, A., Bautista, M.A., Paczan, N., Webb, R., Susskind, J.M.: Hypersim: A photorealistic synthetic dataset for holistic indoor scene understanding. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 10912–10922 (2021)
39. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 10684–10695 (2022)
40. Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton, E.L., Ghasemipour, K., Gontijo Lopes, R., Karagol Ayan, B., Salimans, T., et al.: Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems* **35**, 36479–36494 (2022)
41. Sargent, K., Li, Z., Shah, T., Herrmann, C., Yu, H.X., Zhang, Y., Chan, E.R., Lagun, D., Fei-Fei, L., Sun, D., et al.: Zeronvs: Zero-shot 360-degree view synthesis from a single real image. arXiv preprint arXiv:2310.17994 (2023)
42. Schönberger, J.L., Frahm, J.M.: Structure-from-motion revisited. In: Conference on Computer Vision and Pattern Recognition (CVPR) (2016)
43. Schönberger, J.L., Price, T., Sattler, T., Frahm, J.M., Pollefeys, M.: A vote-and-verify strategy for fast spatial verification in image retrieval. In: Asian Conference on Computer Vision (ACCV) (2016)
44. Shi, Y., Wang, P., Ye, J., Long, M., Li, K., Yang, X.: Mvdream: Multi-view diffusion for 3d generation. arXiv preprint arXiv:2308.16512 (2023)
45. Sølund, T., Buch, A.G., Krüger, N., Aanæs, H.: A large scale 3d object recognition dataset. In: 3DV (2016)
46. Tewari, A., Yin, T., Cazenavette, G., Rezhikov, S., Tenenbaum, J., Durand, F., Freeman, B., Sitzmann, V.: Diffusion with forward models: Solving stochastic inverse problems without direct supervision. *Advances in Neural Information Processing Systems* **36** (2024)
47. Tyszkiewicz, M., Fua, P., Trulls, E.: Disk: Learning local features with policy gradient. *Advances in Neural Information Processing Systems* **33** (2020)
48. Verbin, D., Hedman, P., Mildenhall, B., Zickler, T., Barron, J.T., Srinivasan, P.P.: Ref-nerf: Structured view-dependent appearance for neural radiance fields. In: 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 5481–5490. IEEE (2022)
49. Wang, H., Du, X., Li, J., Yeh, R.A., Shakhnarovich, G.: Score jacobian chaining: Lifting pretrained 2d diffusion models for 3d generation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12619–12629 (2023)

50. Wang, J., Rupprecht, C., Novotny, D.: PoseDiffusion: Solving pose estimation via diffusion-aided bundle adjustment (2023)
51. Wang, P., Tan, H., Bi, S., Xu, Y., Luan, F., Sunkavalli, K., Wang, W., Xu, Z., Zhang, K.: Pflrm: Pose-free large reconstruction model for joint pose and shape prediction. arXiv preprint arXiv:2311.12024 (2023)
52. Wang, Q., Zhou, X., Hariharan, B., Snavely, N.: Learning feature descriptors using camera pose supervision. In: Proc. European Conference on Computer Vision (ECCV) (2020)
53. Wang, S., Leroy, V., Cabon, Y., Chidlovskii, B., Revaud, J.: Dust3r: Geometric 3d vision made easy. arXiv preprint arXiv:2312.14132 (2023)
54. Wang, Z., Lu, C., Wang, Y., Bao, F., Li, C., Su, H., Zhu, J.: Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation. *Advances in Neural Information Processing Systems* **36** (2024)
55. Wang, Z., Bovik, A.C., Sheikh, H.R., Simoncelli, E.P.: Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing* **13**(4), 600–612 (2004)
56. Watson, D., Chan, W., Martin-Brualla, R., Ho, J., Tagliasacchi, A., Norouzi, M.: Novel view synthesis with diffusion models. arXiv preprint arXiv:2210.04628 (2022)
57. Weyand, T., Araujo, A., Cao, B., Sim, J.: Google Landmarks Dataset v2 - A Large-Scale Benchmark for Instance-Level Recognition and Retrieval. In: Proc. CVPR (2020)
58. Wu, R., Mildenhall, B., Henzler, P., Park, K., Gao, R., Watson, D., Srinivasan, P.P., Verbin, D., Barron, J.T., Poole, B., et al.: Reconfusion: 3d reconstruction with diffusion priors. arXiv preprint arXiv:2312.02981 (2023)
59. Wu, X., Averbuch-Elor, H., Sun, J., Snavely, N.: Towers of Babel: Combining images, language, and 3D geometry for learning multimodal vision. In: ICCV (2021)
60. Yang, L., Kang, B., Huang, Z., Xu, X., Feng, J., Zhao, H.: Depth anything: Unleashing the power of large-scale unlabeled data. In: CVPR (2024)
61. Yao, Y., Luo, Z., Li, S., Zhang, J., Ren, Y., Zhou, L., Fang, T., Quan, L.: Blendedmvs: A large-scale dataset for generalized multi-view stereo networks. *Computer Vision and Pattern Recognition (CVPR)* (2020)
62. Yeshwanth, C., Liu, Y.C., Nießner, M., Dai, A.: Scannet++: A high-fidelity dataset of 3d indoor scenes. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 12–22 (2023)
63. Yu, A., Ye, V., Tancik, M., Kanazawa, A.: pixelnerf: Neural radiance fields from one or few images. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4578–4587 (2021)
64. Yu, H.X., Duan, H., Hur, J., Sargent, K., Rubinstein, M., Freeman, W.T., Cole, F., Sun, D., Snavely, N., Wu, J., et al.: Wonderjourney: Going from anywhere to everywhere. arXiv preprint arXiv:2312.03884 (2023)
65. Yu, X., Xu, M., Zhang, Y., Liu, H., Ye, C., Wu, Y., Yan, Z., Liang, T., Chen, G., Cui, S., Han, X.: Mvimngnet: A large-scale dataset of multi-view images. In: CVPR (2023)
66. Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable effectiveness of deep features as a perceptual metric. In: CVPR (2018)
67. Zhou, T., Tucker, R., Flynn, J., Fyffe, G., Snavely, N.: Stereo magnification: Learning view synthesis using multiplane images. *ACM Trans. Graph. (Proc. SIGGRAPH)* **37** (2018), <https://arxiv.org/abs/1805.09817>
68. Zhou, T., Tulsiani, S., Sun, W., Malik, J., Efros, A.A.: View synthesis by appearance flow. In: European Conference on Computer Vision (2016)

69. Zhou, Z., Tulsiani, S.: Sparsefusion: Distilling view-conditioned diffusion for 3d reconstruction. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12588–12597 (2023)