

# A Large-Scale Outdoor Multi-modal Dataset and Benchmark for Novel View Synthesis and Implicit Scene Reconstruction

Chongshan Lu<sup>1</sup> Fukun Yin<sup>1,2</sup> Xin Chen<sup>2</sup> Tao Chen<sup>1\*</sup> Gang Yu<sup>2</sup> Jiayuan Fan<sup>1</sup>  
<sup>1</sup>Fudan University <sup>2</sup>Tencent PCG  
<https://ommo.luchongshan.com>

## Abstract

Neural Radiance Fields (NeRF) has achieved impressive results in single object scene reconstruction and novel view synthesis, as demonstrated on many single modality and single object focused indoor scene datasets like DTU [14], BMVS [41], and NeRF Synthetic [24]. However, the study of NeRF on large-scale outdoor scene reconstruction is still limited, as there is no unified outdoor scene dataset for large-scale NeRF evaluation due to expensive data acquisition and calibration costs.

In this work, we propose a large-scale outdoor multi-modal dataset, **OMMO dataset**, containing complex objects and scenes with calibrated images, point clouds and prompt annotations. A new benchmark for several outdoor NeRF-based tasks is established, such as novel view synthesis, surface reconstruction, and multi-modal NeRF. To create the dataset, we capture and collect a large number of real fly-view videos and select high-quality and high-resolution clips from them. Then we design a quality review module to refine images, remove low-quality frames and fail-to-calibrate scenes through a learning-based automatic evaluation plus manual review. Finally, a number of volunteers are employed to add the text descriptions for each scene and keyframe. Compared with existing NeRF datasets, our dataset contains abundant real-world urban and natural scenes with various scales, camera trajectories, and lighting conditions. Experiments show that our dataset can benchmark most state-of-the-art NeRF methods on different tasks. We will release the dataset and model weights soon.

## 1. Introduction

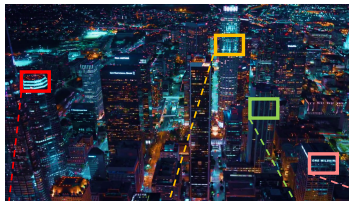
Recent advances in implicit neural representations have achieved remarkable results in photo-realistic novel view synthesis and high-fidelity surface reconstruction [43, 42]. Unfortunately, most of the existing methods focus on a single object or an indoor scene [43, 42, 9, 15, 5], and their

\*Corresponding author.

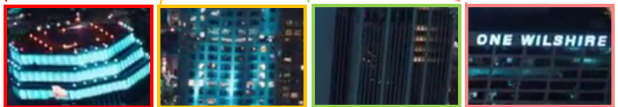
### Prompt annotation

*Buildings at night in the background; two tall buildings in the middle; roads among buildings; cars are passing on the road; tall buildings decorated with blue and orange lights.*

### Calibrated image



### Camera trajectory



### More views



### More scenes



Figure 1. A city scene example from our dataset captured with low illuminance and circle-shaped camera trajectory. We show multi-view calibrated images, the camera track, and text descriptions of the scene. Some details in colored boxes are zoomed in to indicate that our dataset can provide real-world high-fidelity texture details.

synthesis performance will decrease drastically if migrated to outdoor scenes. Although some very recent methods try to solve this problem and are well-designed for large scenes [35, 39], their performance is difficult to compare due to the lack of large-scale outdoor scene datasets and uniform benchmarks.

At present, the existing outdoor scene datasets are either collected with simple scenes containing very few objects,

Table 1. Comparison with existing NeRF datasets, especially those outdoor datasets (or outdoor parts) related to ours. The first group is for single objects, the second group is for large scenes, and the last row is our dataset. For each dataset, we show the number of scenes and images, whether the scene types, camera trajectories, and lighting conditions are diverse, whether they are real-world scenes (called Real), and whether they have multi-modal data (called M-modal).

Datasets	# Scenes	# Images	Types	Camera	Lighting	Real	M-modal
DTU [14]	124	4.2K	No	No	Yes	Yes	No
NeRF [24]	18	3551	No	Yes	No	Yes	No
Scannet [7]	1.5K	2.5M	No	Yes	No	Yes	No
T & T [18]	6	88k	No	No	No	Yes	No
BlendedMVS [41]	28	5k	Yes	No	No	No	No
UrbanScene3D [22]	16	10.4K	Yes	No	No	Part	No
Quad 6k [6]	1	5.1K	No	Yes	No	Yes	No
Mill 19 [35]	2	3.6K	Yes	No	No	Yes	No
<b>Ours</b>	<b>33</b>	<b>14.7K</b>	<b>Yes</b>	<b>Yes</b>	<b>Yes</b>	<b>Yes</b>	<b>Yes</b>

or rendered from virtual scenes, all at a small geographical scale. For example, Tanks and temples [18] provides a benchmark of realistic outdoor scenes captured by a high-precision industrial laser scanner, but its scene scale is still too small ( $463m^2$  on average) and only focuses on a single outdoor object or building. The BlendedMVS [41] and UrbanScene3D [22] datasets contain scene images rendered from reconstructed or virtual scenes, which deviate from the real scene in both texture and appearance details. Collecting images from the Internet can theoretically build very effective datasets [13, 1], like ImageNet [8] and COCO [21], but these methods are not suitable for NeRF-based task evaluation due to the changes of objects and lighting conditions in the scene at different times. Our dataset acquisition method is similar to Mega-NeRF [35], which captures large real-world scenes by drones. But Mega-NeRF only provides two monotonic scenes, which hinders it from being a widely used baseline. Therefore, to our knowledge, no uniform and widely recognized large-scale scene dataset is built for NeRF benchmarking, causing large-scale NeRF research for outdoor far fall behind that for single objects or indoor scenes [14, 41, 24, 7].

To address the lack of large-scale real-world outdoor scene datasets, we introduce a well-selected fly-view multi-modal dataset. The dataset contains totally 33 scenes with prompt annotations, tags, and 14K calibrated images, as shown in Figure 1. Different from the existing methods mentioned above, the sources of our scenes are very extensive, including those collected on the Internet and captured by ourselves. Meanwhile, the collection indicators are also comprehensive and representative, including various scene types, scene scales, camera trajectories, lighting conditions, and multi-modal data that are not available in existing datasets (see Table 1. More importantly, we provide a generic pipeline to generate real-world NeRF-based

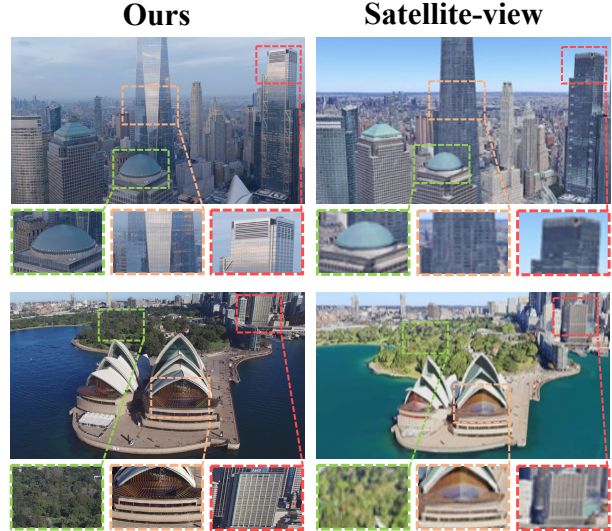


Figure 2. Visual comparison with existing large-scale satellite-view outdoor datasets [39] acquired from Google Earth Studio. The top row is from [39], and the bottom row is corresponding scenes from our fly-view dataset, which is more realistic with clear textures and rich details (zoom-in for the best of views).

data from drone videos on the Internet, which makes our dataset easily to be extensible by the community.

Further, to evaluate the applicability and performance of the built dataset for evaluating mainstream NeRF methods, we build all-around benchmarks including novel view synthesis, scene representations, and multi-modal synthesis based on the dataset. Moreover, we provide several detailed sub-benchmarks for each above task, according to different scene types, scene scales, camera trajectories and lighting conditions, to give a fine-grained evaluation of each method.

To summarize, our main contributions include:

- Aiming at advancing the large-scale NeRF research, we introduce an outdoor scene dataset captured from the real world with multi-modal data, which surpasses all existing relative outdoor datasets in both quantity and diversity, see Table 1 and Sec. 3.3.
- To form a uniform benchmarking standard for outdoor NeRF methods, we create multiple benchmark tasks for mainstream outdoor NeRF methods. Extensive experiments show that our dataset can well support common NeRF-based tasks and provide prompt annotations for future research, see Sec. 4.
- We provide a cost-effective pipeline for converting videos that can be flexibly accessed from the Internet to NeRF-purpose training data, which makes our dataset easily scalable, see Sec. 3.1 and Sec. 3.2.

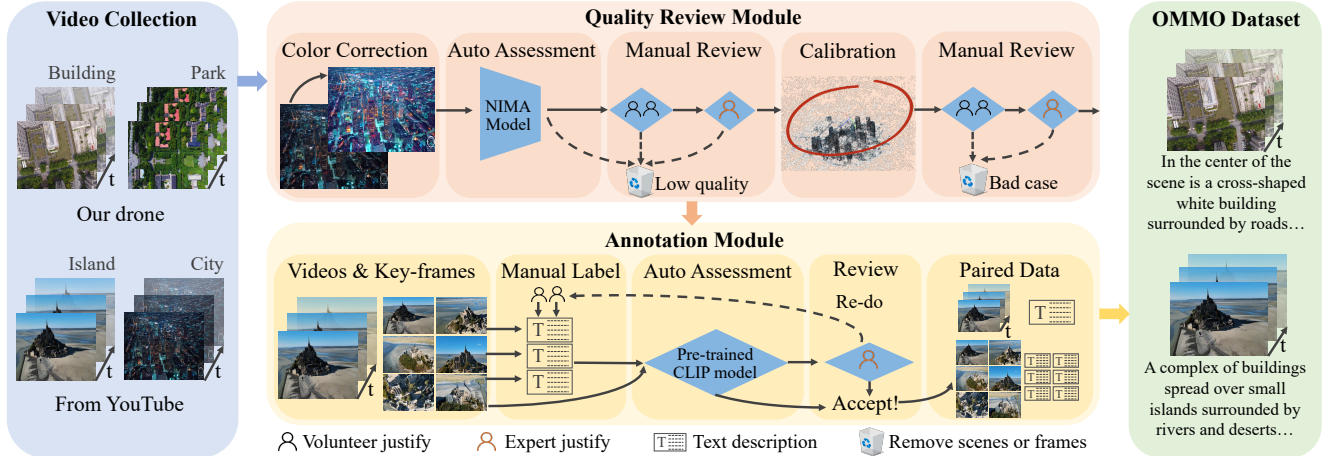


Figure 3. The pipeline for our dataset generation. The original videos are collected from both YouTube and captured by us, and then fed into the review and annotation module. The former mainly removes low-quality frames and failed scenes; the latter annotates text descriptions for scenes and keyframes. Currently, we have generated 33 scenes with 14K images and text descriptions.

## 2. Related Work

### 2.1. Neural Scene Representation and Rendering

Neural Radiance Fields (NeRF) [24] propose an effective implicit neural scene representation method to synthesize novel views by the single scene optimization. Many subsequent coordinate-based methods are inspired by it, and we can classify them into neural surface fields [27, 42, 37] and neural radiance fields [26, 16, 44] based on the shape representation differences. However, both ways have poor performance on large scenes due to the limited representation capability of MLP-based networks adopted in NeRF. Even with prior information, it is difficult to directly apply small-scale scene-focused methods due to the increased scene complexity [43, 19].

Fortunately, some very recent works have started to study the neural representation for large-scale scenes. Mega-NeRF [35] divides the large fly-view scene into multiple small blocks to train specialized NeRFs in parallel. Block-NeRF [34] also adopts this simplified idea, dividing the neighborhood into blocks, and then novel views are sampled from overlapping blocks and combined according to inverse distance weights. CityNerf (BungeeNeRF) [39] introduces a progressive neural radiance field that starts from fitting distant views with a shallow base block, and appends new blocks to accommodate details. NeRF in the Wild (NeRF-W) [23] introduces a series of extensions to NeRF [24] to synthesize novel views of complex scenes, using only unstructured collections of in-the-wild photos. Recursive-NeRF [40] provides an efficient and adaptive rendering and training approach for NeRF, that forwards high uncertainties coordinates to a bigger network with more powerful representational capability.

However, the aforementioned large-scale NeRF methods [35, 34, 39, 23, 40] use different outdoor datasets with various capturing conditions and research focuses, causing difficulty to fairly compare these NeRF methods’ performance on common tasks based on a uniform benchmark.

### 2.2. NeRF-based Datasets and Benchmarks

There are widely used NeRF datasets and established benchmarks for single objects [14, 24], unbounded objects [41], human faces [29], and indoor scenes [7], see the first group in Table 1.

For large-scale outdoor scenes, some datasets provide high-fidelity models from accurate radar scans, but this expensive data acquisition makes the scale and size of these datasets still unsatisfactory [18]. Rendering images from optimized models will result in higher cost and unrealistic scenes [41, 22]. Correspondingly, a low-cost and easy-to-expand way is to collect images of the same scene shot by different people and devices, at different times from the Internet [13, 1]. However, these methods do not meet the needs of common NeRF tasks due to the changes in weather, lighting, and objects in the scene. Mega-NeRF [35] builds two high-quality fly-view real scenes and calibrated images, but has not become a widely used benchmark due to its small size and single type. There are also large-scale NeRF datasets and benchmarks for specific problems, such as neighborhoods [34] or remote sensing [39], see the second group in Table 1.

In conclusion, for the reasons mentioned above, none of the above datasets have formed a widely used uniform benchmark. So a more comprehensive outdoor NeRF-based dataset is required, to facilitate the research and exploitation of larger-scale implicit scene representation. In con-

trast, the built dataset in this work provides 33 large-scale scenes, more than 14K fly-view images with camera poses, rich content, and text descriptions. Meanwhile, several new large-scene fly-view-based benchmarks for novel view synthesis, implicit scene representations, and multi-modal synthesis tasks are also proposed.

### 3. Dataset Generation

Our dataset acquisition, calibration, and annotation pipeline are shown in Figure 3. We first decompose and enhance the original videos by time-sampling and color-correcting, and review the quality of the frames and scene calibrations by automated models followed by volunteers to remove low-quality frames and fail-to-calibrate scenes, see Sec. 3.1. Then volunteers provide prompt annotations for each scene and keyframe, and the CLIP [28] model is exploited to cooperate with human experts, to supervise the semantic consistency between labeled text and corresponding images, see Sec. 3.2. Finally, we introduce the dataset distribution from several aspects such as scene categories and collection cost, see Sec.3.3 and Sec.3.4.

#### 3.1. Acquisition and Calibration Method

**Original Videos.** Our outdoor fly-view dataset mainly comes from two sources: captured by ourselves and collected on YouTube. Among them, the YouTube videos are from worldwide, including urban and natural scenes of various geographical scales, various camera moving trajectories and lighting conditions. However, these videos often have limited resolution due to the compression when uploading, while our captured videos are all 4K HD and can meet more high-fidelity NeRF needs. Benefiting from the diversity of the videos, the proposed scenes in our dataset include various elements, such as buildings, roads, trees, islands, mountains, rivers, etc. In addition, compared to synthetic data, our scenes’ space layout, surface reflection, and lighting conditions are completely real, which supports NeRF methods’ learning about real-world outdoor scenes. To have different lighting conditions in the same scene like DTU [14], we use drones to capture videos of the same building at different times and weather conditions. To explore the impact of different camera trajectories, we scan the same scene with different flight strategies. Overall, our original videos come from 268 real scans from around the world, covering a variety of scene types, camera trajectories, and lighting conditions.

**Color Correction.** We set different sampling intervals according to the video length and frames per second, so that the remaining frame number is between 800-1000 to meet the needs of calibration, yielding a total of 240K frames. For low-light or rainy scenes, we first use an image enhancement model to recover the texture features [17].

**Auto Assessment.** Then NIMA [33] model is employed to evaluate image quality and remove blur, ghosting, and low-quality images. After this step, the 240K frames from 268 scenes in the above step are left with 152K frames.

**Manual Quality Review.** After the quality auto-evaluation, there are still some frames with low quality i.e. blurring, artifacts, or focusing outside the scene. So we employ three volunteers to manually discard frames that do not meet requirements. Specifically, the three volunteers consist of 2 professional data labelers and a domain expert. For each frame, 2 professional data labelers judge whether to keep it or not, and the decision is made if both labelers agree, otherwise it is up to the domain expert for the decision. After this step, there are still 110K frames left.

**Calibration.** Large-scale scene calibration based on highly dynamic images has always been a difficult problem, especially for outdoor scenes with no obvious local detail differences [4, 20, 12]. As a common solution, we use COLMAP [31, 32] to achieve multi-view 3D reconstruction, image calibration and depth image rendering. It is foreseeable that the reconstruction of some scenes with insufficient overlap and textures, or forwardly moving camera motion will fail. These fail-to-calibrate scenes cannot meet the requirements of NeRF-based methods, which need to be removed manually.

**Manual Scene Review.** We invite the above three volunteers to review the calibration quality based on the completeness of scene point clouds, following the same decision-making process. A large number of scenes are discarded at this step, and finally, 33 real-world scans with nearly 14K calibrated images constitute our outdoor dataset, which surpasses existing large-scale datasets in both quantity and diversity.

#### 3.2. Prompt Annotation Method

To provide prompt annotations for multi-modal NeRF, we add text descriptions to each scene and keyframes. Note that generating descriptions is a more subjective and time-consuming task, so we employ more volunteers and apply pre-trained CLIP [28] models.

**Manual Label.** Specifically, six trained volunteers participate in this progress, who non-repeatedly extract and record the corresponding descriptive texts from scenes and keyframes, respectively. These volunteers include high-year Ph.D. students and professors in computer vision and natural language processing fields who can handle this job well.

**Auto Assessment.** Each frame and annotated text is fed into a CLIP [28] model pre-trained on large scene-text datasets [30, 11] to compute their similarity scores. If the similarity score is above the threshold, we accept this annotation. Otherwise, we leave it to experts to double-check.

**Expert Review.** For annotations that the CLIP model

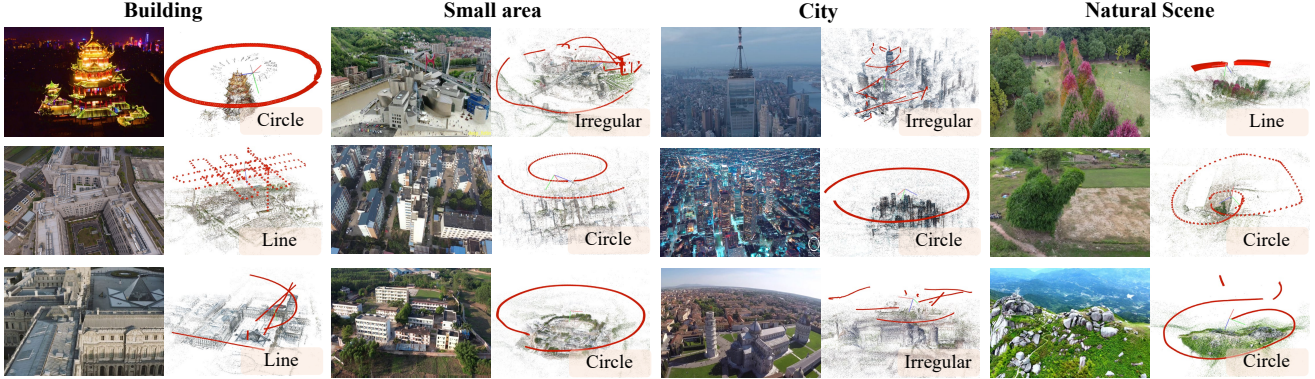


Figure 4. Examples of different types from our dataset. We visualize some scenes and camera trajectories from our dataset, which contain both urban and natural scenes with various scales, camera trajectories, and lighting conditions.

Table 2. Distribution of our dataset. We divide our dataset into subsets based on scene type, camera trajectory, and lighting condition and count the number of each subset.

Scene Type	# Scenes	Camera Trajectory	# Scenes	Lighting Condition	# Scenes
Building	8	Circle	15	Day	30
Small area	9	Line	10	Night	3
City	8	Irregular	8		
Natural scene	8				

cannot judge, experts will evaluate whether the text can describe the image comprehensively and objectively. If possible, we accept the label, otherwise, we hand it over to another volunteer to label again. Fortunately, the frames and scenes in our dataset are only re-annotated at most 2 times.

We have annotated 33 scans with corresponding descriptions and tags, and part of keyframes, which can well meet the training needs of multi-modal NeRF. We are still annotating the remaining keyframes for more potentially complex needs.

### 3.3. Distribution

According to different division methods, the distribution of our dataset is shown in Table 2. However, some scene types are relatively ambiguous. For example, when a building is surrounded by plenty of trees, warehouses, etc., it is difficult to say whether such an image belongs to a building type or not. To resolve the ambiguity, we design a questionnaire and invite 50 voters to determine the attributes of the scenes. Invited volunteers range from 19 to 53 years old, and we recommend a very typical reference for each attribute.

Our dataset contains both natural and urban scenes, which are further divided into buildings, small areas, and whole cities. The performance on different subsets can ver-

ify the most suitable scene type and scale for each NeRF method. Some methods only target scenes with camera trajectories moving in a ring or matrix. For a fair baseline, we divide the scene into circles, lines, and irregular, according to the camera trajectories. In addition, few recent methods focus on low-illuminance NeRF research, so we also provide scenes with different lighting conditions. In particular, we collect some scans of the same scene under different lighting conditions, which can evaluate the ability of the method against poor lighting.

We provide benchmarks for novel view synthesis, generalization, scene representations, and multi-modal synthesis tasks. At the same time, according to the above scene types, our dataset can also produce corresponding sub-benchmarks to evaluate methods under different conditions and settings, see Sec. 4.

### 3.4. Cost

Since our dataset requires volunteers to review images and scenes' quality and annotate the scenes and keyframes with texts, it is unavoidably time-consuming and labor-intensive.

Data collection which involves drone purchasing and shooting, and computing introduced by pre-training scene-text CLIP model, account for a large part of the total cost. But once the drone and CLIP model is ready, we can just use them to prepare more new scenes and add to the dataset, without involving additional cost.

More importantly, with the increase of labeled data that can help train a better model, we can replace part of volunteer review and annotation work with state-of-the-art trained models to optimize our pipeline. Therefore, although our method may be expensive in the early stage, it has good potential to be automated and save future costs that may be introduced when expanding the dataset.



Figure 5. Qualitative visualization results for novel view synthesis (zoom-in for the best of views) on our OMMO dataset.

## 4. Experiments

### 4.1. Setting

To verify the applicability and performance of the built dataset for evaluating NeRF methods, and meanwhile provide a baseline for NeRF-based tasks, we train and evaluate recent NeRF [24], NeRF++ [44], Mip-NeRF [2], Mip-NeRF 360 [3], Mega-NeRF [35] and Ref-NeRF [36] on our datasets.

**NeRF** [24] presents the first continuous MLP-based neural network to represent the scene, that is able to synthesize semantic-consistent novel views by volume rendering.

**NeRF++** [44] separately models the foreground and background neural representations to address the challenge

of modeling large-scale unbounded scenes.

**Mip-NeRF** [2] reduces aliasing artifacts and improves NeRF’s [24] ability to represent fine details, by rendering anti-aliased conical frustums instead of rays.

**Mip-NeRF 360** [3] uses a non-linear scene parameterization, online distillation, and a distortion-based regularizer, to model and produce realistic synthesized views for unbounded real-world scenes.

**Mega-NeRF** [35] proposes a framework for training large-scale 3D scenes by introducing a sparse structure and geometric clustering algorithm, to partition training pixels into different parallel NeRF submodules.

**Ref-NeRF** [36] improves the quality of appearance and normal in synthesized views of the scene, by a new param-

Table 3. Benchmark for novel view synthesis. We present the performance of six state-of-the-art and representative methods on our dataset.  $\uparrow$  means the higher, the better.

Scene ID	Scene Types	Camera Tracks	Lighting Conditions	NeRF [24]			NeRF++ [44]			Mip-NeRF [2]			Mip-NeRF 360 [3]			Mega-NeRF [35]			Ref-NeRF [39]		
				PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
1	Buildings	Irregular	Day	16.93	0.369	0.744	16.86	0.359	0.780	16.84	0.369	0.793	13.91	0.311	0.771	16.12	0.341	0.782	15.10	0.344	0.755
2	Small area	Circles	Day	15.31	0.442	0.694	14.89	0.471	0.653	15.16	0.396	0.731	15.06	0.438	0.646	15.64	0.467	0.679	15.90	0.490	0.632
3	Cities	Lines	Day	14.38	0.278	0.556	14.64	0.294	0.547	14.56	0.288	0.533	14.25	0.309	0.526	15.21	0.325	0.517	15.44	0.371	0.526
4	Buildings	Circle	Night	25.39	0.859	0.431	27.47	0.898	0.380	21.78	0.758	0.469	27.68	0.943	0.292	23.36	0.855	0.419	27.86	0.905	0.404
5	Small area	Circles	Day	22.26	0.670	0.531	24.32	0.729	0.450	14.98	0.544	0.633	25.76	0.801	0.317	25.78	0.763	0.436	23.54	0.706	0.491
6	Natural scenes	Circles	Day	24.09	0.679	0.504	25.59	0.749	0.396	23.18	0.658	0.529	28.86	0.896	0.211	24.92	0.772	0.393	26.07	0.716	0.459
7	Buildings	Lines	Day	5.36	0.166	0.747	21.93	0.707	0.542	15.57	0.643	0.624	23.05	0.734	0.523	22.33	0.691	0.552	25.79	0.731	0.511
8	Cities	Circle	Day	21.14	0.496	0.594	22.91	0.568	0.509	19.82	0.462	0.638	25.07	0.714	0.354	16.65	0.478	0.431	21.21	0.489	0.606
9	Cities	Lines	Day	14.92	0.344	0.744	14.57	0.341	0.732	14.58	0.338	0.746	15.40	0.303	0.706	17.32	0.491	0.673	20.34	0.432	0.649
10	Cities	Irregular	Day	22.26	0.550	0.626	24.37	0.599	0.578	19.80	0.528	0.643	26.68	0.719	0.420	21.78	0.615	0.558	24.23	0.578	0.597
11	Buildings	Circles	Night	22.36	0.816	0.420	24.61	0.852	0.342	22.81	0.822	0.423	27.06	0.931	0.217	24.37	0.844	0.392	23.81	0.843	0.355
12	Small area	Circles	Day	22.41	0.594	0.533	24.29	0.675	0.447	22.13	0.601	0.526	28.12	0.825	0.274	21.60	0.619	0.493	23.06	0.604	0.524
13	Buildings	Lines	Day	22.27	0.592	0.608	23.52	0.623	0.581	18.90	0.537	0.673	26.63	0.771	0.403	25.50	0.722	0.517	23.29	0.605	0.594
14	Small area	Lines	Day	19.85	0.554	0.569	23.89	0.737	0.417	17.06	0.481	0.655	28.06	0.894	0.224	24.42	0.746	0.411	21.76	0.625	0.508
15	Small area	Circles	Day	20.35	0.527	0.552	21.71	0.612	0.490	19.44	0.489	0.594	28.63	0.888	0.179	22.69	0.665	0.445	20.33	0.497	0.576
16	Natural scenes	Circles	Day	17.86	0.397	0.631	18.75	0.405	0.597	18.49	0.399	0.610	10.01	0.344	0.850	20.26	0.532	0.509	19.64	0.428	0.572
17	Natural scenes	Circles	Day	22.02	0.571	0.610	24.20	0.671	0.461	17.01	0.526	0.696	29.53	0.833	0.247	17.23	0.574	0.529	23.17	0.589	0.529
18	Small area	Lines	Day	26.06	0.754	0.428	25.57	0.730	0.461	24.61	0.732	0.469	28.55	0.855	0.265	24.76	0.733	0.448	22.79	0.674	0.569
19	Small area	Circles	Day	14.20	0.399	0.726	13.86	0.373	0.703	13.84	0.394	0.738	14.72	0.367	0.676	23.81	0.682	0.465	14.34	0.386	0.691
20	Cities	Circles	Day	22.84	0.613	0.499	23.28	0.642	0.475	22.41	0.603	0.519	28.33	0.862	0.228	21.11	0.633	0.490	21.54	0.553	0.574
21	Natural scenes	Circles	Day	22.59	0.514	0.532	21.84	0.473	0.593	22.31	0.513	0.537	25.64	0.747	0.344	21.92	0.506	0.578	21.07	0.436	0.672
22	Buildings	Lines	Day	16.53	0.466	0.733	20.66	0.558	0.575	13.37	0.420	0.776	24.79	0.766	0.362	20.84	0.597	0.527	20.31	0.530	0.615
23	Natural scenes	Lines	Day	18.99	0.405	0.669	19.51	0.417	0.597	18.09	0.389	0.671	21.25	0.514	0.539	20.13	0.438	0.585	19.94	0.409	0.622
24	Natural scenes	Lines	Day	19.32	0.386	0.696	23.14	0.522	0.535	16.89	0.374	0.715	25.86	0.707	0.373	23.87	0.563	0.518	22.17	0.452	0.616
25	Natural scenes	Lines	Day	24.72	0.550	0.528	22.42	0.509	0.613	24.24	0.541	0.542	28.91	0.789	0.306	25.98	0.629	0.457	23.62	0.502	0.598
26	Buildings	Irregular	Day	8.56	0.242	0.564	19.94	0.586	0.513	13.43	0.353	0.688	14.59	0.459	0.626	19.23	0.669	0.467	21.00	0.615	0.489
27	Cities	Irregular	Day	4.54	0.006	0.705	21.25	0.548	0.546	14.82	0.453	0.674	21.26	0.599	0.235	20.59	0.606	0.543	20.82	0.519	0.590
28	Small area	Circles	Day	24.48	0.660	0.479	23.28	0.642	0.475	24.76	0.659	0.406	29.62	0.874	0.240	25.87	0.723	0.442	22.17	0.452	0.616
29	Buildings	Circle	Day	22.98	0.608	0.540	23.17	0.617	0.529	23.01	0.609	0.539	25.51	0.740	0.400	21.57	0.611	0.557	21.11	0.543	0.631
30	Natural scenes	Irregular	Day	20.23	0.522	0.605	23.27	0.639	0.476	18.63	0.461	0.675	26.54	0.837	0.296	24.04	0.686	0.459	21.62	0.535	0.586
31	Cities	Circles	Night	18.97	0.365	0.645	19.05	0.371	0.643	18.91	0.358	0.659	13.08	0.234	0.708	20.93	0.596	0.545	19.18	0.372	0.645
32	Cities	Irregular	Day	17.99	0.582	0.621	18.99	0.605	0.540	11.28	0.424	0.687	17.16	0.566	0.601	21.29	0.702	0.475	18.98	0.595	0.565
33	Cities	Irregular	Day	5.79	0.007	0.745	20.19	0.497	0.597	14.31	0.42	0.755	22.76	0.629	0.457	22.89	0.635	0.478	21.23	0.522	0.578
<b>Mean</b>	-	-	-	18.72	0.484	0.600	21.45	0.576	0.538	18.39	0.501	0.623	23.10	0.672	0.418	21.63	0.621	0.508	21.28	0.546	0.574

eterization and structuring of view-dependent outgoing radiance, as well as a regularizer on normal vectors.

**Implementations.** Since there is no official PyTorch implementation of NeRF [24], we use the widely recognized third-party implementation [10]. But for other methods, we use the official implementation from GitHub.

In our dataset, reviewed posed images are numbered from 0 to  $\#image - 1$  under timing sequence, and test views are evenly sampled according to the view ID. That is we sample one from every eight for testing, and the rest are used as training views (i.e., for testing: 0, 8, 16, 24, ...).

All training hyper-parameters follow the original paper’s settings in our experiments. Each scene is trained on a single Nvidia V100 GPU device for around 6-33 hours, depending on the time complexity of each method, and about 32 V100 GPU devices are used in parallel.

**Evaluation Metrics.** To evaluate the performance of each method, we use three common metrics: Peak Signal-to-Noise Ratio (PSNR), Structural Similarity (SSIM) [38], and LPIPS [45] on novel view synthesis. Higher PSNR and SSIM mean better performance, while a lower LPIPS means better.

## 4.2. Novel View Synthesis

**Benchmark.** To establish a benchmark for the large-scale outdoor novel view synthesis, we comprehensively evaluate and report quantitative performances of the above six state-of-the-art methods in our dataset, see Table 3.

It can be seen that except for the failure of NeRF [24] in 4 scenes (7, 26, 27, and 33), other results show that NeRF

can synthesize reasonable novel views, which means that OMMO dataset can support various NeRF-based methods. NeRF++ [44], Mip-NeRF 360 [3], Mega-NeRF [35], and Ref-NeRF [36] perform well on our dataset with an average PSNR of beyond 20, and can maintain the view consistency of each scene, see Figure 5. Among them, Mip-NeRF 360 [3] can synthesize more realistic detailed texture features for large-scale scenes and its quantitative evaluation is more than 6 points higher than other methods on PSNR, SSIM, LPIPS. Our benchmarks are open to all NeRF-based methods, and we are also ready to evaluate newer large-scale scene NeRF methods once they are proposed.

In particular, we notice that most of the scenes where NeRF fails are based on irregular camera trajectories, which suggests that NeRF may be more suitable for scenes captured with stronger trajectory consistency constraints and more overlap (such as equidistant circular acquisitions). So we divide OMMO dataset into subsets according to the data types, and provide sub-benchmarks to study the most suitable setting for each method.

**Sub-benchmark split by scene types.** According to different scales of urban and natural scenes, we propose 4 sub-benchmarks for buildings, small areas, cities and natural scenes. It can be seen from Table 4 that all methods perform worse in cities than in smaller-scale subsets, i.e. buildings and small areas. These performance differences show that the large-scale scene implicit representation is still not as well resolved as for single objects or small scenes.

**Sub-benchmark split by camera tracks.** Circular camera trajectories tend to present better experiment perfor-

Table 4. More sub-benchmarks for novel view synthesis. We divide our dataset into subsets based on different scene types, camera trajectories, and lighting conditions, and provide sub-benchmarks under different settings.  $\uparrow$  means the higher, the better.

Scene ID	Sub-benchmark	NeRF [24]			NeRF++ [44]			Mip-NeRF [2]			Mip-NeRF 360 [3]			Mega-NeRF [35]			Ref-NeRF [39]			
		PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	
1,4,7,8,11,13,22,26	<b>Buildings</b>	17.32	0.501	0.605	22.24	0.644	0.528	17.82	0.546	0.636	22.85	0.704	0.444	21.05	0.650	0.511	22.30	0.633	0.541	
2,5,12,14,15,18,19,28,29		<b>Small areas</b>	20.88	0.579	0.561	21.66	0.621	0.514	19.44	0.545	0.588	24.89	0.742	0.358	22.90	0.668	0.486	20.56	0.553	0.582
3,8,9,10,20,27,31,32,33		<b>Cities</b>	15.87	0.360	0.637	19.92	0.496	0.574	16.72	0.430	0.650	20.44	0.548	0.471	19.75	0.565	0.523	20.33	0.492	0.592
6,16,17,21,23,24,25,30	<b>Natural scenes</b>	21.23	0.503	0.597	22.34	0.548	0.534	19.86	0.483	0.622	24.58	0.708	0.396	22.29	0.588	0.504	22.16	0.508	0.582	
2,4,5,6,8,11,12,15,16,17,19,20,21,28,31	<b>Circles</b>	21.08	0.573	0.559	22.00	0.609	0.508	19.80	0.545	0.581	23.81	0.713	0.386	21.74	0.647	0.483	21.53	0.564	0.556	
3,7,9,13,14,18,22,23,24,25	<b>Lines</b>	18.24	0.450	0.628	20.99	0.544	0.560	17.79	0.474	0.640	23.68	0.664	0.423	22.04	0.594	0.521	21.55	0.533	0.581	
1,10,26,27,29,30,32,33	<b>Irregular</b>	14.91	0.361	0.644	21.01	0.556	0.570	16.52	0.452	0.682	21.05	0.608	0.476	20.94	0.608	0.540	20.51	0.531	0.599	
ALL-{4, 11, 31}	<b>Day</b>	18.37	0.465	0.610	21.23	0.563	0.547	18.12	0.487	0.634	23.15	0.670	0.420	21.51	0.607	0.514	21.05	0.531	0.585	
4, 11, 31	<b>Night</b>	22.24	0.680	0.499	23.71	0.707	0.455	21.17	0.646	0.517	22.61	0.703	0.406	22.89	0.765	0.452	23.62	0.707	0.468	



Figure 6. Different scans of the same scene during day and night. Both RGB and depth images are synthesized by Mega-NeRF [35].

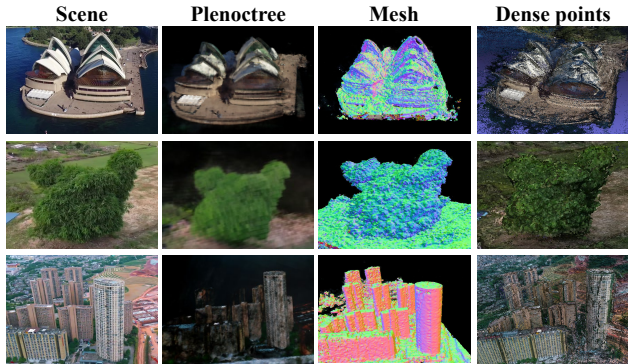


Figure 7. Examples of various scene representations from our dataset through different methods.

mance than other types, especially irregular ones. This is because 360-degree views contain richer texture features from different angles, and the focus of views is overlapped to maintain the view consistency.

**Sub-benchmark split by lighting conditions.** Intuitively, daytime scenes are richer in texture and easier to learn their representation than dark ones. However, we find that almost every method performs better on the low-light subset than on the normal-light subset. We visualize two different scans of the same scene generated by Mega-NeRF [35] during day and night, as shown in Figure 6. It is not difficult to see that in low light settings, the implicit network uses black areas to erase details when generating the RGB images, which reduces the synthesis difficulty and

tricks the evaluation metrics, while synthesised poor depth map illustrates the network’s incapacity to understand and represent the scene. So efficient low-light NeRF methods are urged to solve this problem.

### 4.3. Scene Representation

To evaluate the performance of our dataset on surface or scene reconstruction tasks, we reconstruct scenes with different representations by using a variety of methods including implicit networks. Specifically, plenocree, mesh, and dense points are provided by Mega-NeRF [35], Instant-NGP [25], and Colmap [31, 32], respectively. It can be seen from Figure 7 that neither the implicit network nor the feature matching reconstruction method can reconstruct the large scene finely. Theoretically, the advantage of implicit scene representation is that, the scene can be reconstructed with high resolution benefiting from the continuous representation. So the scene representation benchmarks of large-scale outdoor scenes based on NeRF are still to be build. Please refer to the supplementary materials for more results.

## 5. Discussion

**Conclusion.** We introduce a well-selected large-scale outdoor multi-modal fly-view dataset, OMMO, to address the problem of no widely-used benchmark for outdoor NeRF-based methods. The built OMMO surpasses the previous datasets in several key indicators such as quantity, quality and variety, by providing 33 real-world scenes with more than 14K posed images and text description. With the help of our cost-effective data collection pipeline, it is easy to expand our dataset by continuously converting new internet videos into NeRF-purpose training data. We provide benchmarks on multiple tasks such as novel view synthesis, implicit scene representations, and multi-modal synthesis by evaluating existing methods. Experiments show that our dataset can well support mainstream NeRF-based tasks.

**Limitation.** Scenes with low-illumination, rain and fog are still few in the current dataset due to the limited calibration and reconstruction ability of COLMAP [31, 32]. We will try more reconstruction and calibration methods to solve this problem. Meanwhile, we are continuing to expand our dataset thanks to our cost-effective pipeline.



## References

- [1] Sameer Agarwal, Yasutaka Furukawa, Noah Snavely, Ian Simon, Brian Curless, Steven M Seitz, and Richard Szeliski. Building rome in a day. *Communications of the ACM*, 54(10):105–112, 2011.
- [2] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5855–5864, 2021.
- [3] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5470–5479, 2022.
- [4] Aljaz Bozic, Pablo Palafox, Justus Thies, Angela Dai, and Matthias Nießner. Transformerfusion: Monocular rgb scene reconstruction using transformers. *Advances in Neural Information Processing Systems*, 34:1403–1414, 2021.
- [5] Di Chen, Yu Liu, Lianghua Huang, Bin Wang, and Pan Pan. Geoaug: Data augmentation for few-shot nerf with geometry constraints. In *European Conference on Computer Vision*, pages 322–337. Springer, 2022.
- [6] David Crandall, Andrew Owens, Noah Snavely, and Dan Huttenlocher. Discrete-continuous optimization for large-scale structure from motion. In *CVPR 2011*, pages 3001–3008. IEEE, 2011.
- [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.
- [8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [9] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised nerf: Fewer views and faster training for free. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12882–12891, 2022.
- [10] Github. nerf-pytorch. <https://github.com/yenchenlin/nerf-pytorch>.
- [11] Michael Grubinger, Paul Clough, Henning Müller, and Thomas Deselaers. The iapr tc-12 benchmark: A new evaluation resource for visual information systems. In *International workshop on image*, volume 2, 2006.
- [12] Haoyu Guo, Sida Peng, Haotong Lin, Qianqian Wang, Guofeng Zhang, Hujun Bao, and Xiaowei Zhou. Neural 3d scene reconstruction with the manhattan-world assumption. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5511–5520, 2022.
- [13] Jared Heinly, Johannes L Schonberger, Enrique Dunn, and Jan-Michael Frahm. Reconstructing the world\* in six days\*(as captured by the yahoo 100 million image dataset). In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3287–3295, 2015.
- [14] Rasmus Jensen, Anders Dahl, George Vogiatzis, Engin Tola, and Henrik Aanæs. Large scale multi-view stereopsis evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 406–413, 2014.
- [15] Mohammad Mahdi Johari, Yann Lepoittevin, and François Fleuret. Geonerf: Generalizing nerf with geometry priors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18365–18375, 2022.
- [16] Kacper Kania, Kwang Moo Yi, Marek Kowalski, Tomasz Trzcziński, and Andrea Tagliasacchi. Conerf: Controllable neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18623–18632, 2022.
- [17] Hanul Kim, Su-Min Choi, Chang-Su Kim, and Yeong Jun Koh. Representative color transform for image enhancement. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4459–4468, 2021.
- [18] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Transactions on Graphics (ToG)*, 36(4):1–13, 2017.
- [19] Jonáš Kulháněk, Erik Derner, Torsten Sattler, and Robert Babuška. Viewformer: Nerf-free neural rendering from few images using transformers. *arXiv preprint arXiv:2203.10157*, 2022.
- [20] Quentin Legros, Julián Tachella, Rachael Tobin, Aongus McCarthy, Sylvain Meignen, Gerald S Buller, Yoann Altmann, Stephen McLaughlin, and Michael E Davies. Robust 3d reconstruction of dynamic scenes from single-photon lidar using beta-divergences. *IEEE Transactions on Image Processing*, 30:1716–1727, 2020.
- [21] 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 *European conference on computer vision*, pages 740–755. Springer, 2014.
- [22] Yilin Liu, Fuyou Xue, and Hui Huang. Urbanscene3d: A large scale urban scene dataset and simulator. *arXiv preprint arXiv:2107.04286*, 2021.
- [23] Ricardo Martin-Brualla, Noha Radwan, Mehdi SM Sajjadi, Jonathan T Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7210–7219, 2021.
- [24] 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. In *ECCV*, 2020.
- [25] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Trans. Graph.*, 41(4):102:1–102:15, July 2022.
- [26] Michael Niemeyer, Jonathan T Barron, Ben Mildenhall, Mehdi SM Sajjadi, Andreas Geiger, and Noha Radwan. Regnerf: Regularizing neural radiance fields for view synthesis

- from sparse inputs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5480–5490, 2022.
- [27] Michael Oechsle, Songyou Peng, and Andreas Geiger. Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5589–5599, 2021.
- [28] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021.
- [29] Eduard Ramon, Gil Triginer, Janna Escur, Albert Pumarola, Jaime Garcia, Xavier Giro-i Nieto, and Francesc Moreno-Noguer. H3d-net: Few-shot high-fidelity 3d head reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5620–5629, 2021.
- [30] Ramon Sanabria, Ozan Caglayan, Shruti Palaskar, Desmond Elliott, Loïc Barrault, Lucia Specia, and Florian Metzger. How2: a large-scale dataset for multimodal language understanding. In *Proceedings of the Workshop on Visually Grounded Interaction and Language (ViGIL)*. NeurIPS, 2018.
- [31] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4104–4113, 2016.
- [32] Johannes L Schönberger, Enliang Zheng, Jan-Michael Frahm, and Marc Pollefeys. Pixelwise view selection for unstructured multi-view stereo. In *European conference on computer vision*, pages 501–518. Springer, 2016.
- [33] Hossein Talebi and Peyman Milanfar. Nima: Neural image assessment. *IEEE transactions on image processing*, 27(8):3998–4011, 2018.
- [34] Matthew Tancik, Vincent Casser, Xinchun Yan, Sabeek Pradhan, Ben Mildenhall, Pratul P Srinivasan, Jonathan T Barron, and Henrik Kretzschmar. Block-nerf: Scalable large scene neural view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8248–8258, 2022.
- [35] Haithem Turki, Deva Ramanan, and Mahadev Satyanarayanan. Mega-nerf: Scalable construction of large-scale nerfs for virtual fly-throughs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12922–12931, 2022.
- [36] Dor Verbin, Peter Hedman, Ben Mildenhall, Todd Zickler, Jonathan T Barron, and Pratul P Srinivasan. Ref-nerf: Structured view-dependent appearance for neural radiance fields. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5481–5490. IEEE, 2022.
- [37] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *Advances in Neural Information Processing Systems*, 34:27171–27183, 2021.
- [38] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004.
- [39] Yuanbo Xiangli, Linning Xu, Xingang Pan, Nanxuan Zhao, Anyi Rao, Christian Theobalt, Bo Dai, and Dahua Lin. Bungeenerf: Progressive neural radiance field for extreme multi-scale scene rendering. In *The European Conference on Computer Vision (ECCV)*, volume 2, 2022.
- [40] Guo-Wei Yang, Wen-Yang Zhou, Hao-Yang Peng, Dun Liang, Tai-Jiang Mu, and Shi-Min Hu. Recursive-nerf: An efficient and dynamically growing nerf. *IEEE Transactions on Visualization and Computer Graphics*, 2022.
- [41] Yao Yao, Zixin Luo, Shiwei Li, Jingyang Zhang, Yufan Ren, Lei Zhou, Tian Fang, and Long Quan. Blendedmvs: A large-scale dataset for generalized multi-view stereo networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1790–1799, 2020.
- [42] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces. *Advances in Neural Information Processing Systems*, 34:4805–4815, 2021.
- [43] Fukun Yin, Wen Liu, Zilong Huang, Pei Cheng, Tao Chen, and Gang YU. Coordinates are not lonely—codebook prior helps implicit neural 3d representations. *arXiv preprint arXiv:2210.11170*, 2022.
- [44] Kai Zhang, Gernot Riegler, Noah Snively, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*, 2020.
- [45] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018.

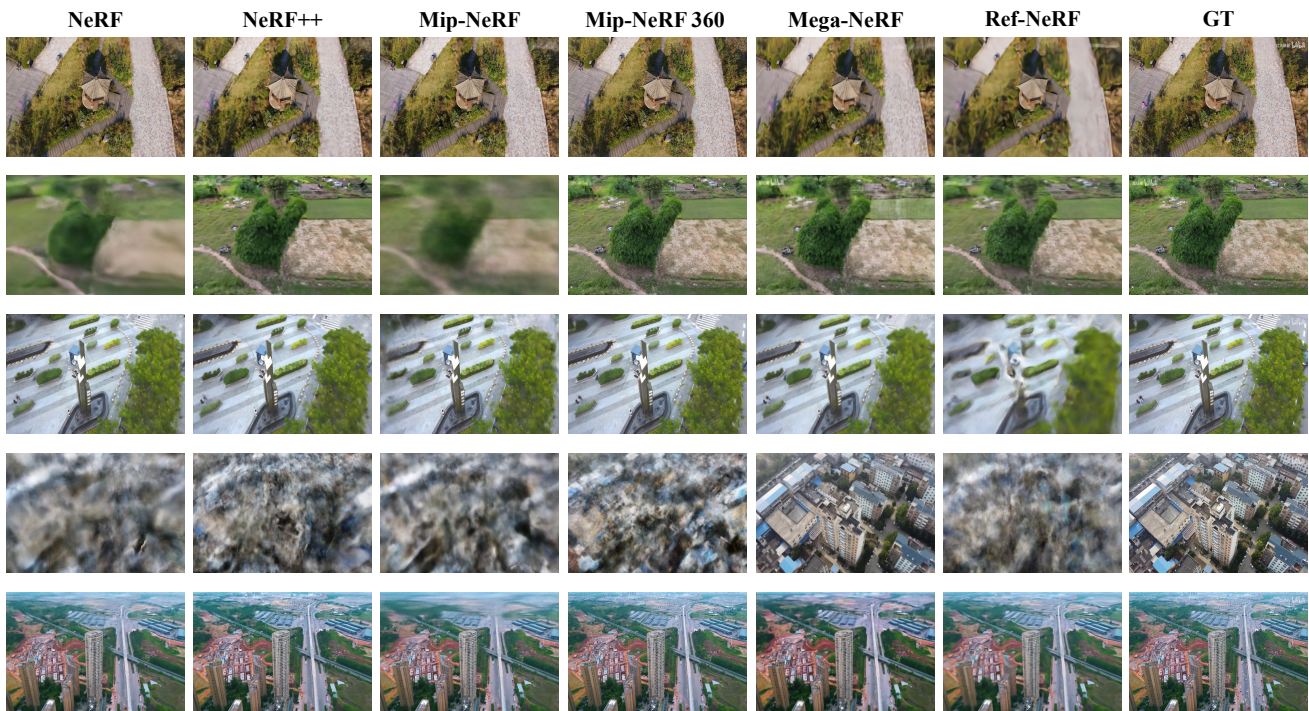
# Appendix

In this supplementary material, we provide the appendix section and a supplemental video to better understand our database and benchmarks. This appendix involves more qualitative or experimental results (Sec. B), details of our dataset generation method (Sec. C), and dataset analysis (Sec. D). The supplemental video contains a brief introduction to our dataset, some examples in detail, and more comprehensive synthesis results in surrounding views or progressive views.

## A. More Qualitative Results

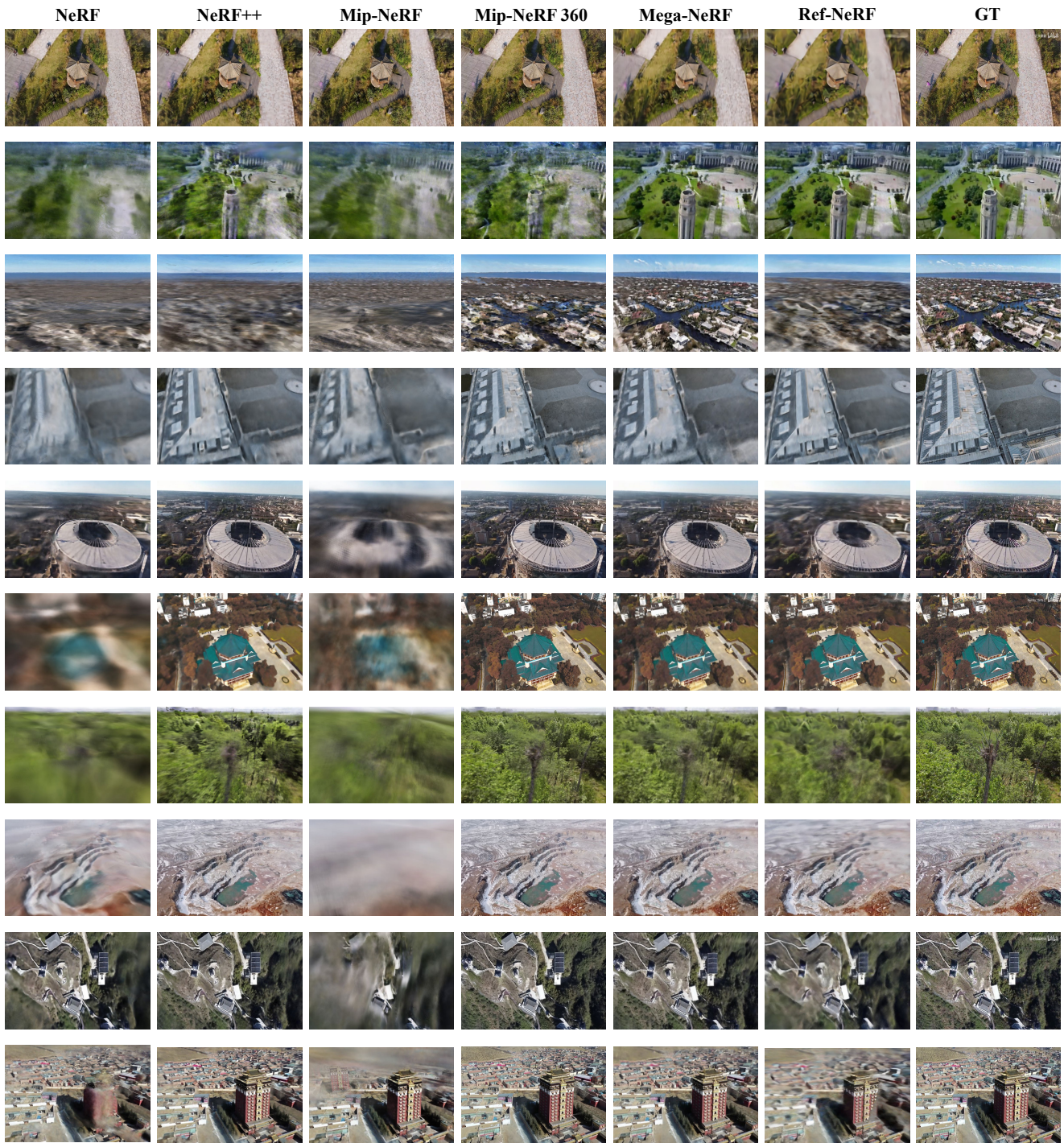
### A.1. Novel View Synthesis

In our main manuscript, we can only provide the visualization results of five scenes due to the length limitation. In this section, more qualitative results are presented to demonstrate the novel view synthesis ability of each method, see Fig. 1.



Part 1 / 2

Figure 1. More qualitative visualization results for novel view synthesis (zoom-in for the best of views) on our OMMO dataset.



Part 2/2

Figure 1. More qualitative visualization results for novel view synthesis (zoom-in for the best of views) on our OMMO dataset.

## A.2. Scene Representation

To further demonstrate that our OMMO dataset can well support surface or scene reconstruction tasks including NeRF-based methods, we visualize more shape results by various representations in Fig. 2. Among them, plenoctree, mesh, and dense points are provided by Mega-NeRF [35], InstantNGP [25], and Colmap [31, 32], respectively.

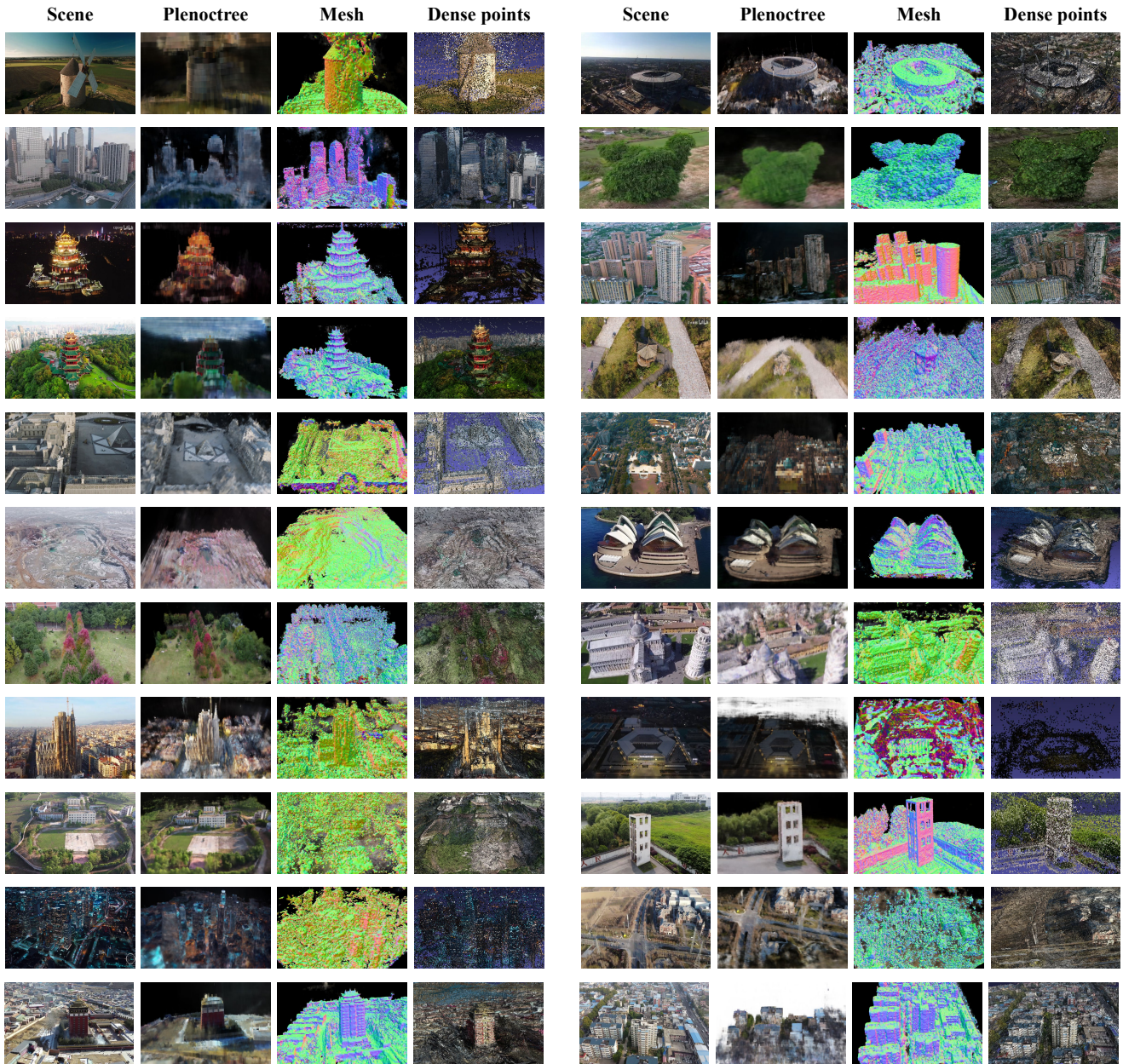


Figure 2. More qualitative visualization results for various scene representations (zoom-in for the best of views) through the state-of-the-art methods on the OMMO dataset.

## B. More Experiments

**Multi-modal NeRF Synthesis.** Since there is no available NeRF-based method for text-assisted fidelity novel view synthesis, inspired by CoCo-INR [43], we replace its image-based pre-scene codebook with text-based codebook and apply it in NeRF [24] and CoCo-INR [43] as our benchmark. Specifically, we apply a text-based attentional coordinate module in front of the last MLP layer of volume rendering network, where the text-based codebook is encoded by our textual prompts through the pre-trained CLIP [28] model, see Fig. 3.

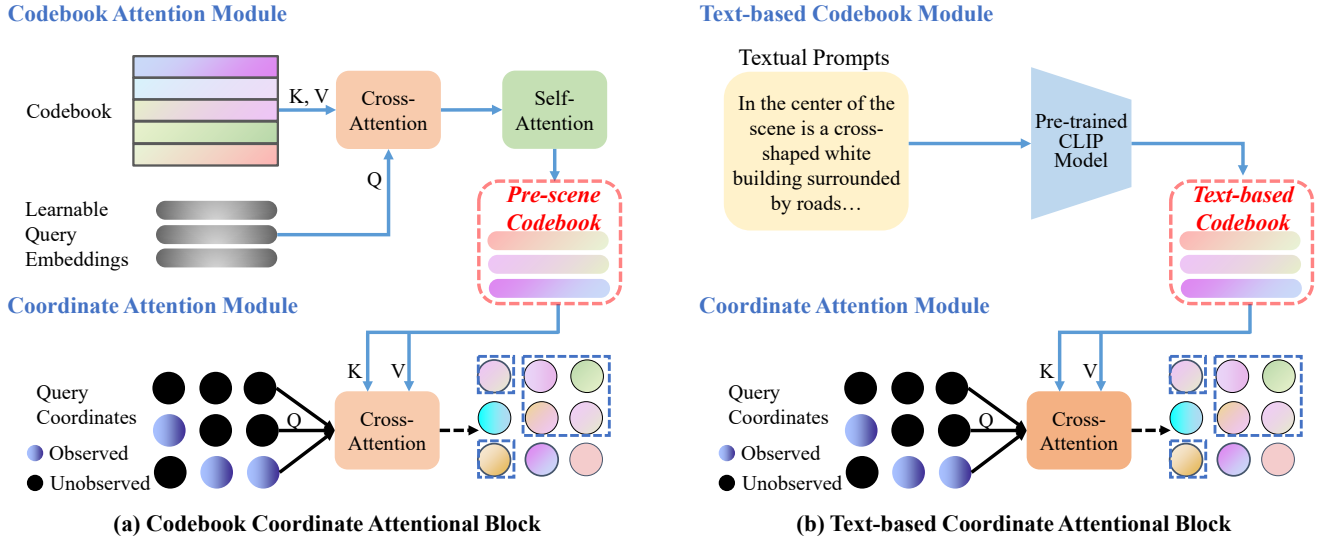


Figure 3. Structure of coordinate attentional blocks. The left sub-figure (a) is the CoCo-INR’s [43] codebook coordinate attentional block, which extracts image features related to the current scene from the prior by codebook attention module to form a pre-scene codebook. The right sub-figure (b) is our text-based coordinate attentional block, which obtains the scene-related text-based codebook by encoding textual prompts of each scene. Both inject scene-related features into each coordinate through the coordinate attention module.

Table 1. Performance comparison of with or without textual prompts for novel view synthesis on our OMMO dataset. We report the results of each method and the performance improvement after injecting textual prompts.  $\uparrow$  means the higher, the better.

Mtehod	Without Textual Prompts			With Textual Prompts			Improvement (%) $\uparrow$		
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR	SSIM	LPIPS
NeRF [24]	18.72	0.484	0.600	19.01	0.500	0.591	1.5	3.2	1.5
CoCo-INR [43]	16.80	0.489	0.681	16.97	0.490	0.678	1.0	0.2	0.4

It can be seen from Table 1, even without a well-designed module for injecting text information, the performance of both NeRF [24] and CoCo-INR [43] methods have improved. Since the textual prompts contain more global features about rich geometry or appearance information, which are shared by different views in the scene to guarantee the network to generate view-consistency results. We hope to inspire more image-text multi-modal NeRF methods to synthesize photo-realistic rendering results and decent geometry by exploring effective ways to make full use of textual prompts. The benchmark on each scene and the sub-benchmarks on different scene types are shown in Table 2 and Table 3.

Table 2. Benchmark for multi-modal NeRF synthesis. We present the performance of text-assisted novel view synthesis based on existing methods on our OMMO dataset.  $\uparrow$  means the higher, the better.

Scene ID	Scene Types	Camera Tracks	Lighting Conditions	NeRF [24] w/o Prompts			NeRF [24] w/ Prompts			CoCo-INR [43] w/o Prompts			CoCo-INR [43] w/ Prompts		
				PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
1	Buildings	Irregular	Day	16.93	0.369	0.744	16.89	0.366	0.729	14.31	0.432	0.788	14.81	0.431	0.785
2	Small area	Circles	Day	15.31	0.442	0.694	15.61	0.465	0.711	16.04	0.597	0.632	16.25	0.597	0.626
3	Citys	Lines	Day	14.38	0.278	0.556	14.42	0.277	0.573	15.59	0.485	0.616	16.62	0.509	0.585
4	Buildings	Circle	Night	25.39	0.859	0.431	24.94	0.851	0.425	21.61	0.876	0.480	21.87	0.879	0.481
5	Small area	Circles	Day	22.26	0.670	0.531	21.31	0.652	0.564	18.16	0.657	0.597	20.08	0.675	0.573
6	Natural scenes	Circles	Day	24.09	0.679	0.504	23.78	0.655	0.535	19.65	0.630	0.576	19.39	0.627	0.578
7	Buildings	Lines	Day	5.36	0.166	0.747	6.25	0.183	0.697	16.53	0.628	0.679	15.38	0.567	0.654
8	Citys	Circle	Day	21.14	0.496	0.594	21.55	0.510	0.571	16.94	0.413	0.687	16.57	0.407	0.704
9	Citys	Lines	Day	14.92	0.344	0.744	15.02	0.345	0.749	13.70	0.340	0.773	13.68	0.340	0.765
10	Citys	Irregular	Day	22.26	0.550	0.626	22.44	0.551	0.624	18.81	0.536	0.694	18.62	0.535	0.693
11	Buildings	Circles	Night	22.36	0.816	0.420	22.58	0.820	0.412	17.08	0.746	0.494	17.35	0.747	0.491
12	Small area	Circles	Day	22.41	0.594	0.533	22.80	0.608	0.512	17.87	0.475	0.658	17.81	0.473	0.659
13	Buildings	Lines	Day	22.27	0.592	0.608	23.12	0.619	0.576	16.55	0.532	0.698	17.02	0.542	0.671
14	Small area	Lines	Day	19.85	0.554	0.569	20.73	0.591	0.534	15.44	0.485	0.663	15.19	0.482	0.665
15	Small area	Circles	Day	20.35	0.527	0.552	20.70	0.549	0.533	16.37	0.407	0.702	16.45	0.411	0.689
16	Natural scenes	Circles	Day	17.86	0.397	0.631	17.53	0.362	0.647	15.37	0.384	0.633	15.24	0.376	0.640
17	Natural scenes	Circles	Day	22.02	0.571	0.610	22.23	0.575	0.596	20.52	0.575	0.619	19.38	0.527	0.648
18	Small area	Lines	Day	26.06	0.754	0.428	26.48	0.770	0.402	17.31	0.527	0.658	17.35	0.532	0.664
19	Small area	Circles	Day	14.20	0.399	0.726	14.19	0.397	0.720	15.41	0.388	0.701	15.82	0.413	0.694
20	Citys	Circles	Day	22.84	0.613	0.499	23.30	0.636	0.465	18.28	0.434	0.676	18.09	0.431	0.685
21	Natural scenes	Circles	Day	22.59	0.514	0.532	22.99	0.541	0.508	17.08	0.358	0.744	17.28	0.359	0.720
22	Buildings	Lines	Day	16.53	0.466	0.733	20.404	0.539	0.598	14.86	0.408	0.759	14.73	0.406	0.772
23	Natural scenes	Lines	Day	18.99	0.405	0.669	19.09	0.405	0.671	17.57	0.335	0.673	17.43	0.332	0.701
24	Natural scenes	Lines	Day	19.32	0.386	0.696	18.52	0.379	0.708	18.63	0.347	0.765	18.27	0.341	0.814
25	Natural scenes	Lines	Day	24.72	0.550	0.528	25.24	0.576	0.496	20.15	0.434	0.717	20.29	0.434	0.711
26	Buildings	Irregular	Day	8.56	0.242	0.564	8.56	0.242	0.564	9.19	0.336	0.924	9.23	0.341	0.913
27	Citys	Irregular	Day	4.54	0.006	0.705	4.91	0.249	0.818	16.19	0.443	0.699	16.07	0.443	0.687
28	Small area	Circles	Day	24.48	0.660	0.479	24.32	0.630	0.493	20.12	0.536	0.643	21.13	0.595	0.621
29	Buildings	Circle	Day	22.98	0.608	0.540	23.58	0.631	0.516	16.57	0.439	0.733	17.93	0.453	0.716
30	Natural scenes	Irregular	Day	20.23	0.522	0.605	21.02	0.559	0.569	12.36	0.431	0.760	15.40	0.450	0.719
31	Citys	Circles	Night	18.97	0.365	0.645	19.09	0.371	0.634	17.88	0.465	0.685	17.57	0.459	0.704
32	Citys	Irregular	Day	17.99	0.582	0.621	18.00	0.582	0.628	17.01	0.623	0.588	16.94	0.622	0.590
33	Citys	Irregular	Day	5.79	0.007	0.745	5.79	0.007	0.744	15.20	0.436	0.761	14.68	0.431	0.770
<b>Mean</b>	-	-	-	<b>18.72</b>	<b>0.484</b>	<b>0.600</b>	<b>19.01</b>	<b>0.500</b>	<b>0.592</b>	<b>16.80</b>	<b>0.489</b>	<b>0.681</b>	<b>16.97</b>	<b>0.490</b>	<b>0.678</b>

Table 3. More sub-benchmarks for multi-modal NeRF synthesis. We divide our dataset into subsets based on different scene types, camera trajectories, and lighting conditions, and provide sub-benchmarks under different settings.  $\uparrow$  means the higher, the better.

Scene ID	Sub-benchmark	NeRF [24] w/o Prompts			NeRF [24] w/ Prompts			CoCo-INR [43] w/o Prompts			CoCo-INR [43] w/ Prompts		
		PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
1,4,7,8,11,13,22,26	<b>Buildings</b>	17.32	0.501	0.605	18.04	0.516	0.572	15.88	0.546	0.689	15.87	0.540	0.684
2,5,12,14,15,18,19,28,29	<b>Small areas</b>	20.88	0.579	0.561	21.08	0.588	0.554	17.03	0.501	0.665	17.56	0.515	0.656
3,8,9,10,20,27,31,32,33	<b>Citys</b>	15.87	0.360	0.637	16.06	0.392	0.645	16.62	0.464	0.687	16.54	0.464	0.687
6,16,17,21,23,24,25,30	<b>Natural scenes</b>	21.23	0.503	0.597	21.30	0.507	0.591	17.67	0.437	0.686	17.84	0.431	0.691
2,4,5,6,8,11,12,15,16,17,19,20,21,28,31	<b>Circles</b>	21.08	0.573	0.559	21.13	0.575	0.555	17.89	0.529	0.635	18.02	0.532	0.634
3,7,9,13,14,18,22,23,24,25	<b>Lines</b>	18.24	0.450	0.628	18.93	0.468	0.600	16.63	0.452	0.700	16.60	0.449	0.700
1,10,26,27,29,30,32,33	<b>Irregular</b>	14.91	0.361	0.644	15.15	0.398	0.649	14.96	0.460	0.743	15.46	0.463	0.734
ALL-{4, 11,31}	<b>Day</b>	18.37	0.465	0.610	18.69	0.482	0.602	16.59	0.468	0.694	16.77	0.469	0.690
4, 11,31	<b>Night</b>	22.24	0.680	0.499	22.20	0.681	0.490	18.86	0.696	0.553	18.93	0.695	0.559

### C. Method Details

We show some dropped frames or scenes during dataset generation to better understand our selection and review standard in Fig. 4. At auto assessment stage, the image quality assessment model [33] is employed to remove frames with blur, artifacts, ghosting and incorrect colors caused by overexposure or optical effects. In this way, about 64% of the frames remained, but there are still some low-quality frames with blur, subtitles, abnormal brightness or transparency caused by fading in or out at the beginning or end of the video. So during the manual quality review process, volunteers and experts will work together to remove these frames. After scene calibration and reconstruction, some scenes will fail, such as with insufficient overlap and textures, or forwardly moving camera motion. These fail-to-calibrate scenes cannot meet the requirements of NeRF-based methods, which need to be removed at the manual scene review stage.

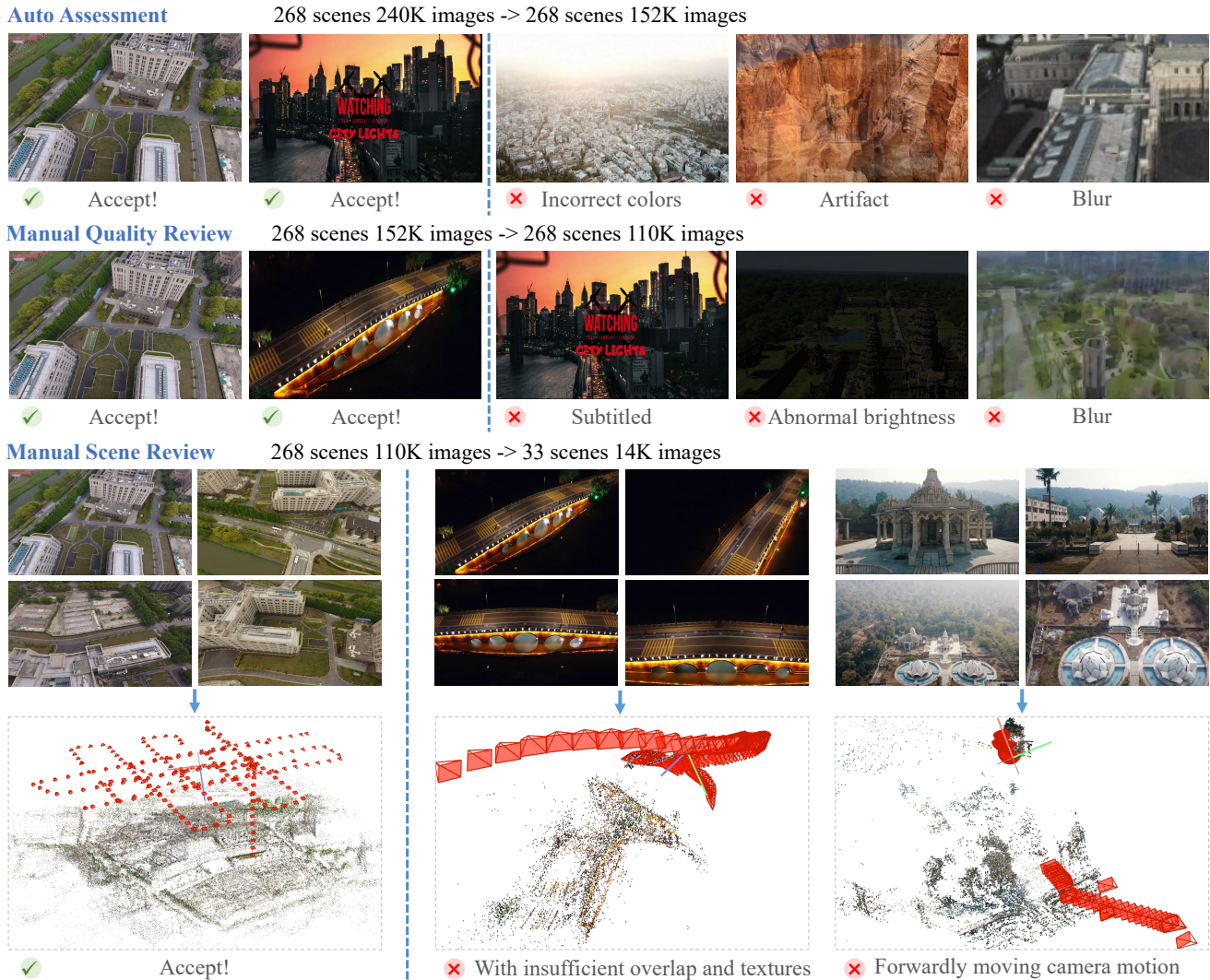


Figure 4. Some examples of dropped frames or scenes at auto assessment, manual quality review, and manual scene review stages. Meanwhile, we also show the number of images and scenes before and after the review at each stage.

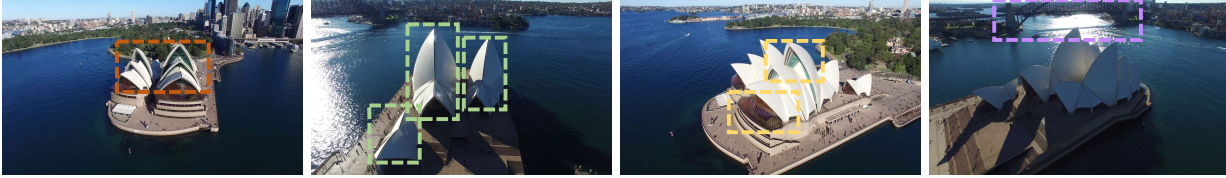


## D. Dataset Analysis

### D.1. Textual Prompts

We show an example of scene prompt annotations from our OMMO dataset in Fig. 5. Our prompts annotation comprehensively describes every detail of the scene center and its surrounding environment in many short sentences.

#### Views



#### Prompt annotation

- The **white building** is captured by a circular camera track.
- The building is located on a peninsula surrounded by water on three sides.
- The shape of the building is **three shell-shaped sub-buildings**, two of which are juxtaposed with larger shells and another one is smaller.
- The two larger sub-buildings are composed of four pointed shells in a cascade.
- The smaller sub-building consists of two back-to-back shells.
- The **glass between each layer of shells is yellow or green**.
- There are many people around the building.
- There is a white carport in front of the building, where some cars are parked;
- Behind the building is a round island with many trees planted on it.
- There is **a bridge across the river** on one side of the building.
- There are dense buildings on both sides of the bridge.

Figure 5. **Textual Prompt.** An example of annotations. Several phrases and their corresponding patches are highlighted in the same color.

We report the word statistic for all scene prompts annotations (only including nouns that appear more than 4 times), as shown in Fig. 6. It can be seen that our data distribution is comprehensive and reasonable, including building, buildings (architectural complex), trees, roads, lawn and rivers, etc. Meanwhile, the number of keywords can roughly reflect the distribution of different scenes, such as natural scene: urban scene (building, small area, city) is about 1:3.

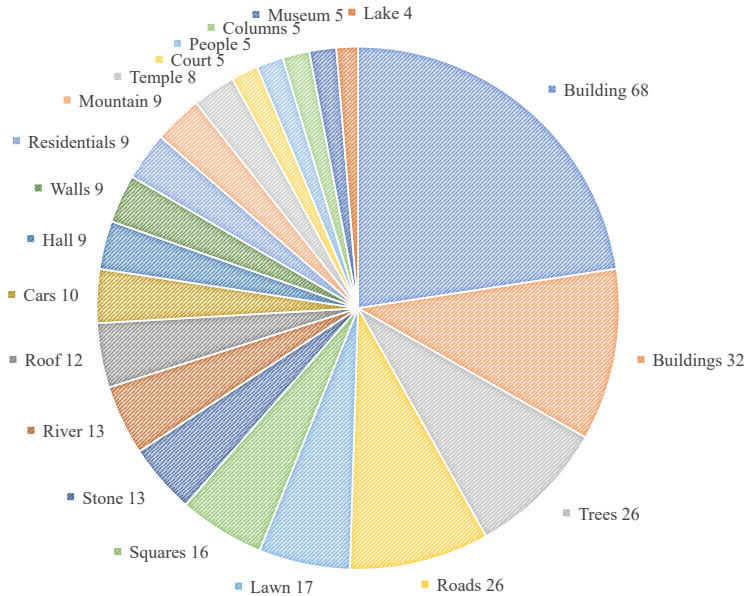


Figure 6. **Word statistic.** Only include nouns that appear more than 4 times in our OMMO dataset.

## D.2. User Instructions

Our OMMO dataset structure list is shown below. The first-level directory contains the scene list, the training and validation split file, and sub-folders for each scene. Each scene contains original video and sub-folders for images, camera matrices and textual prompts.

