# XScale-NVS: Cross-Scale Novel View Synthesis with Hash Featurized Manifold

Guangyu Wang[1], Jinzhi Zhang[1], Fan Wang[2], Ruqi Huang[1,†], Lu Fang[1,†]

Tsinghua University[1], Alibaba Group[2]

(a) Meta Representation [39]   (b) 3D Gaussian Splatting [16]   (c) ZipNeRF [5]   (d) **Ours (Micro-scale)**   (e) **Ours (Macro-scale)**
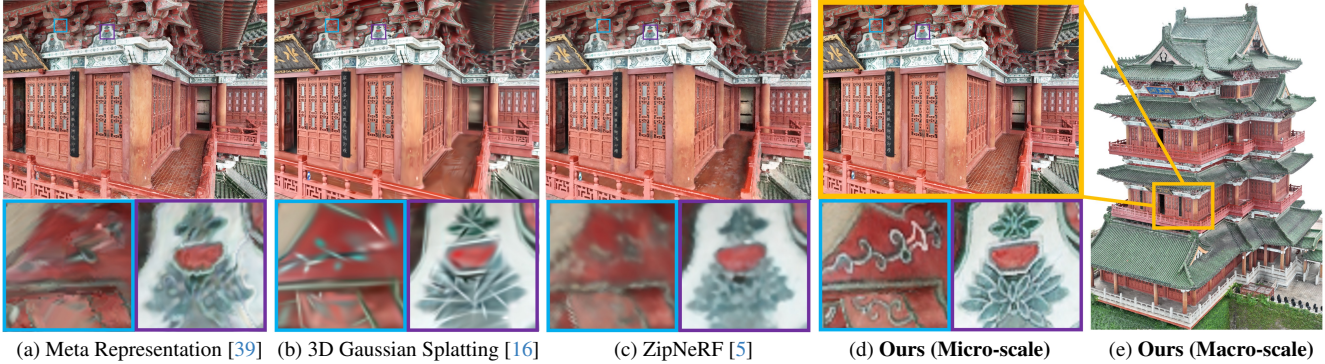
Figure 1. We propose hash featurized manifold representation for high-fidelity cross-scale neural rendering of real-world large-scale scenes. Compared to the recent advances [5, 16, 39], our method synthesizes novel views with unprecedented levels of realism. Please zoom in to see the high-quality details.

## Abstract

*We propose* XScale-NVS *for high-fidelity cross-scale novel view synthesis of real-world large-scale scenes. Existing representations based on explicit surface suffer from discretization resolution or UV distortion, while implicit volumetric representations lack scalability for large scenes due to the dispersed weight distribution and surface ambiguity. In light of the above challenges, we introduce hash featurized manifold, a novel hash-based featurization coupled with a deferred neural rendering framework. This approach fully unlocks the expressivity of the representation by explicitly concentrating the hash entries on the 2D manifold, thus effectively representing highly detailed contents independent of the discretization resolution. We also introduce a novel dataset, namely* GigaNVS, *to benchmark cross-scale, high-resolution novel view synthesis of real-world large-scale scenes. Our method significantly outperforms competing baselines on various real-world scenes, yielding an average LPIPS that is ∼ 40% lower than prior state-of-the-art on the challenging* GigaNVS *benchmark. Please see our project page at: xscalenvs.github.io.*

## 1. Introduction

Neural rendering has attracted significant amount of research interests and influenced down-stream applications including virtual reality, visual effects, robotic simulation, to name a few. Among the recent advances, a notable trend focuses on free-view rendering of real-world large-scale scenes [13, 23, 29, 35, 37, 39]. While these approaches excel at recovering macro-scale structures from images captured at a distance, they often struggle to simultaneously deliver micro-scale details, leading to a discrepancy that hinders comprehensive visual perception of the real world (See Fig. 1). This limitation necessitates a novel scene representation tailored for high-fidelity cross-scale novel view synthesis (NVS) of real-world large-scale scenes. In particular, we identify and address two critical challenges towards this goal as follows.

First of all, reconstructing real-world large-scale scenes in high quality requires to collect images from both distant and close-up views. However, current real-world NVS benchmarks [4, 6, 17, 33, 51] are limited to a single macro-scale setup that neglects close-up imagery for local details and focus primarily on small-scale scenes. To this end, we introduce GigaNVS dataset to benchmark cross-scale, high-resolution NVS of real-world large-scale scenes. The dataset contains seven scenes covering an average area of $1.4 \times 10^6 m^2$, each of which is captured using a combi-

nation of aerial and ground photography, yielding a collection of 1,600 ∼ 18,000 high-quality 5K/8K multi-view images per scene with unstructured scale variations. Remarkably, our dataset captures millimeter-level details from scenes with square-kilometer-level areas, enabling an effective texture resolution of 30 billion during the reconstruction. Therefore, the proposed `GigaNVS` benchmark simultaneously characterizes gigantic scene scale and unparalleled richness of real-world contents by providing cross-scale, high-resolution, and real-captured data.

Secondly, there lacks a suitable representation simultaneously possessing global robustness and local expressivity. Formulations that favor the former typically rely on an imperfect surface proxy reconstructed from multi-view stereo (MVS) [10, 15, 32, 46], and then featurize the proxy on either the parametrized 2D UV map [23, 36, 39] or the 3D surface [2, 16, 18, 27, 31, 45, 53]. While being global-structure-aware, such representations typically struggle to represent the local intricate details. For instance, as shown in Fig. 1 (a), UV-based featurizations [23, 36, 39] suffer from the inherent distortion of surface parametrization [8, 14, 28], which severely degrades the rendering quality. Meanwhile, representations based on more global 3D discretization, such as point cloud [2, 18, 27, 31, 53], mesh [45], or Gaussians [16], all suffer from limited featuremetric resolutions, i.e., the resolutions of the local spatial features, which depend on the respective discretization resolution. As shown in Fig. 1 (b), such feature allocation mode can not fully describe the micro details presented in real-world large-scale scenes.

On the other end of the spectrum, state-of-the-art implicit volumetric representations [5, 21, 25, 30, 42] offer the potential to express arbitrarily high spatial resolutions with multi-resolution hash encoding [25], enabling more effective and efficient neural field reconstruction. Similar to earlier neural field based methods [24, 29, 35, 37, 44], these approaches rely on volume rendering, which alpha-composites multiple samples along the ray, to produce effective gradient signals for a plausible geometric optimization [40]. Unfortunately, for complex large-scale scenes, this comes at the cost of a dispersed weight distribution and leads to surface ambiguities, since the informative surface intersection is compromised with non-informative, multi-view inconsistent samples, as shown in Fig. 1 (c).

To address the limitations of existing scene representations, we introduce *hash featurized manifold* representation, i.e., an expressive hash-based featurization upon the surface manifold for high-fidelity cross-scale NVS of real-world large scenes. Our key insight is to featurize the 2D surface manifold with 3D volumetric hash encoding, sidestepping the complex surface parametrization to strictly preserve geometric conformality while leveraging rasterization to concentrate the learnable hash entries on multi-view

consistent signals throughout the optimization. Compared to existing representations based on explicit 3D discretization, our hash-based featurization unlocks the dependence on the discretization resolution by adaptively erasing the impact of unimportant spatial features on hash entries, and naturally bypasses any surface parametrizations to circumvent the distortions. In turn, by explicitly prioritizing the sparse manifold instead of densely featurizing a redundant volume, the expressivity of multi-resolution hash encoding is substantially incentivized. The reason is that the extensively correlated spatial features receive clean colour gradients purely from multi-view consistent regions without any disturbance from inconsistent regions, thus better reasoning about the optimal capacity for modelling view-invariant reflectance components.

Hash featurized manifold can be tightly coupled with a deferred neural rendering pipeline to simultaneously advance expressivity and efficiency. To render our representation, we first obtain a screen-space feature buffer through rasterization, and then employ an Multi-Layer Perceptrons (MLPs) based neural shader to reason about the view dependent surface colour. We also introduce two enhancements tailored for our representation to better express the cross-scale rich contents: 1) A surface multisampling scheme to enable a prefiltered featurization, which copes with the unstructured scale variations and eliminates aliasing by sampling the curved surface; 2) A manifold deformation mechanism to implicitly enforce multi-view consistency for a better tolerance of geometric imperfections on the surface manifold. Combining the two designs together, our representation essentially expresses a deformable frustum near the surface rather than a single ray-surface intersection, allowing for more flexible descriptions of fine-grained details.

The proposed method significantly outperforms prior approaches on the challenging `GigaNVS` benchmark and the public Tanks&Temples dataset [17]. Remarkably, our method reduces the average LPIPS relative to the state-of-the-art by 40% on `GigaNVS`, pushing the boundary of in-the-wild cross-scale neural rendering towards unprecedented levels of details and realism. In summary, our main contributions are as follows:

- We propose hash featurized manifold representation that fully unleashes the expressivity of volumetric hash encoding by rasterizing the surface manifold to explicitly prioritize multi-view consistency.
- We present a deferred neural rendering framework to efficiently decode the representation and propose two tailored designs to better describe cross-scale details.
- We introduce `GigaNVS` dataset to benchmark cross-scale, high-resolution NVS of real-world large-scale scenes, where our method demonstrates significant improvements over prior approaches.

## 2. Related Work

**Large-scale Scene NVS.** Recent approaches exploit spatial partitioning [13, 35, 37] and geometric priors [20, 23, 29, 39] to better handle large-scale scenes. However, these works can only represent scenes reasonably at a macro-scale yet exhibit excessively blurry artifacts when navigating closer to inspect micro details. Remarkably, BungeeNeRF [44] enables multi-scale neural rendering of large scenes with a progressive optimization scheme to gradually expand the model and training data. However, it requires an explicit split of scales among the input images, thus being limited to remote sensing like scenarios, where the scale can be readily measured by the flight altitude. By contrast, we focus on general large-scale scenes and unstructured scale variations commonly presented in practical perception. We hold that our task can not be well addressed without modifying the fundamental scene representation.

**Large-scale Scene Benchmark.** The recent trend of 3D datasets [19, 22, 26, 44, 47] start to focus on large scenes, however, the available multi-view images are either synthesized by game engines [19, 22, 26] or re-rendered from reconstructed mesh [22, 44, 47], drastically deteriorating the diversity and fidelity of real-world contents. Empowered by the recent advance of gigapixel-level sensation [41, 49, 50], the GigaMVS dataset [51] captures real-world large-scale scenes with ultra-high-resolution imagery. However, the collected images are shot from a distance, whose amount is also limited due to the complicated imaging procedure. To the best of our knowledge, the proposed `GigaNVS` is the first real-captured dataset targeting cross-scale, high-resolution NVS of large-scale scenes.

**Representations upon explicit 3D discretization.** Thies et al. [36] incorporate neural textures into traditional mesh rasterization pipeline and use a CNN-based neural renderer to enable high quality NVS. Another line of approaches [2, 18, 27, 31, 53] follow a similar pipeline but use point as the surface primitive and directly featurize the surface in 3D. Recently, 3D Gaussian Splatting [16] demonstrates great success in terms of rendering quality and efficiency with a highly flexible point-based representation.

**Implicit Volumetric Representations.** The exploding neural representations [3–5, 21, 24, 25, 40, 48] implicitly encode scene geometry as the density field [3–5, 24, 25, 34] or signed distance field [21, 30, 40, 42, 48] and represent appearance as the radiance field. Notably, iNGP [25] proposes multi-resolution hash encoding, which can conceptually represent a dense feature grid at arbitrarily high resolution by hashing a fixed-size learnable array. Neuralangelo [21] further extends this with a coarse-to-fine control of the hash grid and demonstrates impressive neural surface reconstruction quality. However, existing volumetric neural fields can not generalize well to large-scale scenarios due to the inherent surface ambiguities.

## 3. Methodology

In this section, we introduce hash featurized manifold, aiming at high-fidelity cross-scale NVS of real-world large-scale scenes. Given a set of posed multi-view images $\{\mathcal{I}_k\}$ and a mesh $\mathcal{S}$ reconstructed using off-the-shelf MVS techniques, we conduct an expressive and distortion-free surface-based featurization with multi-resolution hash encoding. We then propose a rasterization pipeline and a neural shader tailored to the novel featurization for efficient, robust and high quality rendering.

In the following, we first review the basic building blocks of our method (Sec. 3.1), then formulate our representation (Sec. 3.2) and introduce several specific designs (Sec. 3.3) to further enhance the rendering quality.

### 3.1. Preliminaries

**Deferred Neural Rendering.** Approaches like [2, 16, 18, 23, 27, 31, 36, 39, 45, 53] exploit various featurizations upon the explicit 3D discretization of the scene (typically an imperfect MVS reconstruction), and seamlessly incorporate graphics rasterization pipeline to enable photo-realistic neural rendering.

Given an explicit surface proxy $\mathcal{S}$ in the form of the triangle mesh or point cloud, these methods augment $\mathcal{S}$ with UV-based featurization $\mathcal{F}_{uv} : \mathbb{R}^2 \mapsto \mathbb{R}^Z$ or 3D-surface-based featurization $\mathcal{F}_s : \mathbb{R}^3 \mapsto \mathbb{R}^Z$, where learnable $Z$-dimensional feature descriptors are assigned to each surface primitive, which can be the UV texel [23, 36, 39], mesh vertex [45], 3D point [2, 18, 27, 31, 53], or Gaussian [16]. Taking 3D-surface-based methods [2, 16, 18, 27, 31, 45, 53] as an example, given a target camera pose for rendering, the rasterizer $\mathcal{R}_s$ assigns each pixel a 3D point $\boldsymbol{x}_s \in \mathcal{S}$ at which the ray traced from the pixel intersects with $\mathcal{S}$, and then efficiently samples the corresponding feature $\mathcal{F}_s(\boldsymbol{x}_s)$. This essentially results in a screen-space feature buffer $\mathcal{R}_s(\{\mathcal{F}_s(\boldsymbol{x}_s)\}) \in \mathbb{R}^{H \times W \times Z}$. Finally, a decoder $\mathcal{M}$ is utilized to interpret the feature buffer to synthesize the final RGB rendering $\mathcal{I} \in \mathbb{R}^{H \times W \times 3}$, which can be parametrized as CNNs [2, 27, 31, 53], MLPs [18], or spherical harmonic (SH) composition [16], optionally taking as input the view direction map $\{\boldsymbol{d}_i\} \in \mathbb{R}^{H \times W \times 3}$. Putting pieces together, the deferred neural rendering framework can be formulated as:

$$\mathcal{I} = \mathcal{M}(\mathcal{R}_s(\{\mathcal{F}_s(\boldsymbol{x}_s)\}), \{\boldsymbol{d}_i\}). \tag{1}$$

**Multi-resolution Hash Encoding.** State-of-the-art neural field based methods [5, 21, 25, 30, 42] leverage a hybrid representation, i.e., a combination of hash encoding [25] and a shallow MLP-based decoder, to efficiently reconstruct the geometry and appearance.

In practice, hash encoding is conducted in a multi-resolution manner to effectively handle the collision. Let $\mathcal{V}$ be the continuous volume to be featurized, we denote by

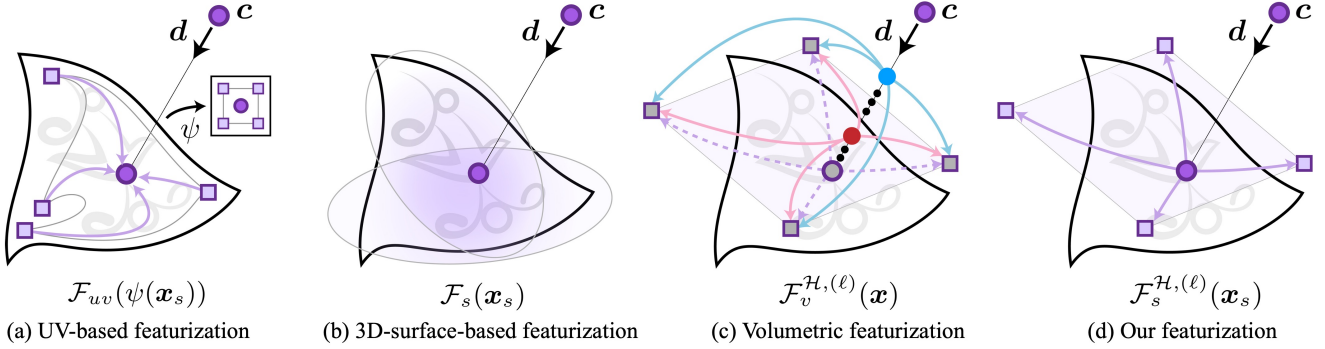| $\mathcal{F}_{uv}(\psi(\boldsymbol{x}_s))$ | $\mathcal{F}_s(\boldsymbol{x}_s)$ | $\mathcal{F}_v^{\mathcal{H},(\ell)}(\boldsymbol{x})$ | $\mathcal{F}_s^{\mathcal{H},(\ell)}(\boldsymbol{x}_s)$ |
|:---:|:---:|:---:|:---:|
| (a) UV-based featurization | (b) 3D-surface-based featurization | (c) Volumetric featurization | (d) Our featurization |

Figure 2. Illustration of different featurizations. The curved triangle represents a micro surface patch in 3D with rich texture details, which can be reflected by the close-up imagery of a real-world large scene. Let the bold purple dot represent the target pixel colour $\boldsymbol{c}$, and we denote by the black arrow $\boldsymbol{d}$ the view direction of the camera ray. (a) UV-based featurizations [23, 36, 39] tend to disorganize the feature distribution due to distortions [8, 14, 28] in surface parametrization $\psi$. (b) Existing 3D-surface-based featurizations [2, 16, 18, 27, 31, 45, 53] fail to express the sub-primitive-scale intricate details given the limited discretization resolution. (c) Volumetric featurizations [5, 21, 25, 30, 42] inevitably yield a dispersed weight distribution during volume rendering, where many multi-view inconsistent yet highly weighted samples ambiguate surface colour and deteriorate surface features with inconsistent colour gradient. (d) Our method leverages hash encoding to unlock the dependence of featuremetric resolution on discretization resolution, and utilizes rasterization to fully unleash the expressivity of volumetric hash encoding by propagating clean and multi-view consistent signals to surface features.

$\mathcal{F}_v^{\mathcal{H},(\ell)} : \mathbb{R}^3 \mapsto \mathbb{R}^Z$ the $\ell$-th level volumetric hash encoding, which conceptually expresses a dense 3D feature grid by hashing a learnable array $\mathcal{H} \in \mathbb{R}^Z$ of fixed length $N_h$. Given a 3D point $\boldsymbol{x} \in \mathcal{V}$, the corresponding hash feature $\mathcal{F}_v^{\mathcal{H},(\ell)}(\boldsymbol{x}) \in \mathbb{R}^Z$ is queried by tri-linearly interpolating the hash entries at the vertices of the grid cell encompassing $\boldsymbol{x}$. The hash features across all resolutions are concatenated as: $\mathcal{F}_v^{\mathcal{H}}(\boldsymbol{x}) = \{\mathcal{F}_v^{\mathcal{H},(\ell)}(\boldsymbol{x})|_{\ell=1}^L\} \in \mathbb{R}^{LZ}$, which is then passed to the light-weight decoder to reason about the implicit field.

## 3.2. Hash Featurized Manifold

Despite the notable advancements, existing scene representations become problematic in capturing the cross-scale, in-the-wild richness of large-scale scenes. In the following, we start by identifying several critical limitations of existing representations, which are illustrated in Fig. 2 for better intuitions.

For UV-based methods (Fig. 2 (a)), the severely distorted parametrization, which is commonly encountered in large-scale scenarios with highly complicated shapes, essentially leads to a disorganized feature distribution on the 3D surface $\mathcal{S}$ without preserving the conformality, thus leading to stretched and blurry artifacts at local details. On the other hand, as shown in Fig. 2 (b), existing representations based on explicit 3D discretization only assign a single feature descriptor to each surface primitive, thus failing to faithfully describe the intricate details within the surface primitive (e.g., the elliptic gaussians shaded in purple). For implicit volumetric representations (Fig. 2 (c)), the weight distribution of volume rendering is scattered throughout the optimization, i.e., there exist many multi-view inconsistent yet highly weighted samples (e.g., the red and blue dots) that

contaminate the supervision of surface colour and mislead the adaptations of surface features by propagating inaccurate colour gradient.

To address the above limitations, we propose a novel scene representation, namely hash featurized manifold, aiming for the exploration of a more expressive surface-based featurization leveraging multi-resolution hash encoding and deferred neural rendering.

We now give the detailed description of our design. Similar to existing neural representations based on explicit 3D discretization, we first reconstruct a mesh using off-the-shelf MVS techniques to serve as a 3D surface proxy of the scene. Then, we compute the bounding volume $\mathcal{V}$ of the mesh $\mathcal{S}$ and featurize it with volumetric multi-resolution hash encoding $\mathcal{F}_v^{\mathcal{H},(\ell)} : \mathbb{R}^3 \mapsto \mathbb{R}^Z$, which gives us a hash featurized volume $\{\mathcal{F}_v^{\mathcal{H},(\ell)}(\boldsymbol{x})\}_{\boldsymbol{x} \in \mathcal{V}}$. Throughout the optimization, we leverage the mesh rasterizer $\mathcal{R}_s$ to calculate the 3D surface intersection $\boldsymbol{x}_s$ for each pixel and query the multi-resolution hash feature $\mathcal{F}_v^{\mathcal{H},(\ell)}(\cdot) \in \mathbb{R}^Z$ only at the surface intersections $\boldsymbol{x}_s \in \mathcal{S}$ instead of in the redundant volume $\boldsymbol{x} \in \mathcal{V}$. With the explicit guidance of $\mathcal{S}$, the learnable hash table $\mathcal{H} \in \mathbb{R}^Z$ is forced to prioritize multi-view consistent surface regions $\{\boldsymbol{x} \in \mathcal{S}|\boldsymbol{x} \in \mathcal{V}\}$ with the most important fine scale features, inherently turning the redundant volumetric featurization into an expressive surface-based featurization:

$$\{\mathcal{F}_v^{\mathcal{H},(\ell)}(\boldsymbol{x})\}_{\boldsymbol{x} \in \mathcal{V}} \mapsto \{\mathcal{F}_s^{\mathcal{H},(\ell)}(\boldsymbol{x}_s)\}_{\boldsymbol{x}_s \in \mathcal{S}}. \quad (2)$$

Compared to vanilla hash-based volumetric featurization $\mathcal{F}_v^{\mathcal{H},(\ell)}$ (Fig. 2 (c)), our featurization $\mathcal{F}_s^{\mathcal{H},(\ell)}$ (Fig. 2 (d)) samples a single surface intersection along the pixel ray, eliminating the surface colour ambiguity. Moreover, our
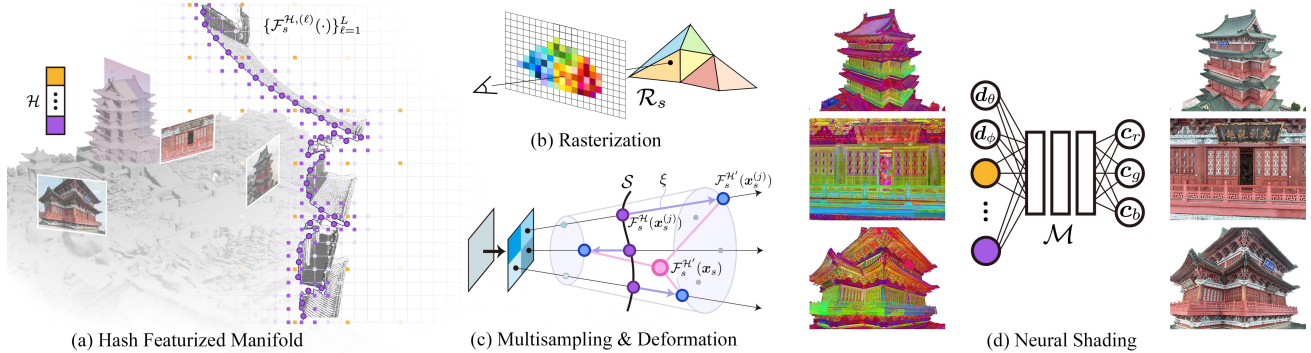
Figure 3. An overview of the hash featurized manifold representation and our neural rendering framework. (a) We first reconstruct the scene as a mesh using MVS and featurize the surface manifold with volumetric multi-resolution hash encoding. (b) We then rasterize the featurized manifold into screen space and (c) optionally perform surface multisampling and manifold deformation to express a deformable frustum for a better representation of the cross-scale details. (d) An MLP-based neural shader decodes the rasterized feature buffer and account for the view dependent colour. Remarkably, we leverage rasterization to concentrate the featurization on multi-view consistency throughout the optimization, inherently converting the redundant volumetric featurization into an expressive surface-based featurization.

representation essentially concentrates on multi-view consistent surface and propagates clean, accurate gradient signals to surface features throughout the optimization, thus boosting the expressivity of multi-resolution hash encoding to faithfully describe the surface colour. Different from existing 3D-surface-based featurizations $\mathcal{F}_s$ (Fig. 2 (b)), our featurization utilizes surface-aware hash encoding to effectively capture the sub-primitive-scale details regardless of the discretization resolution, demonstrating superior scalability towards large-scale scenes and cross-scale contents. In addition, our method allocates hash features on regular 3D voxel grids without relying on surface parametrizations, circumventing the distortion issues in UV-based featurizations $\mathcal{F}_{uv}$ (Fig. 2 (a)).

An overview of our deferred neural rendering pipeline is illustrated in Fig. 3. Let $\mathcal{F}_s^{\mathcal{H}}(\boldsymbol{x}_s) = \left\{ \mathcal{F}_s^{\mathcal{H},(\ell)}(\boldsymbol{x}_s)\big|_{\ell=1}^{L} \right\} \in \mathbb{R}^{LZ}$ be the concatenation of hash features across all encoding resolutions, and given a target camera pose, we start by rasterizing the hash featurized manifold into screen space to create a hash-based feature buffer $\mathcal{R}_s(\{\mathcal{F}_s^{\mathcal{H}}(\boldsymbol{x}_s)\}) \in \mathbb{R}^{H \times W \times LZ}$. To enable realistic NVS, we use a light-weight MLP-based decoder $\mathcal{M} : \mathbb{R}^{LZ} \times \mathbb{R}^3 \mapsto \mathbb{R}^3$ to compute the view-dependent colour, taking as inputs the resulting hash features and the view direction map $\{\boldsymbol{d}_i\}$. The proposed neural rendering pipeline is therefore defined as:

$$\mathcal{I} = \mathcal{M}(\mathcal{R}_s(\{\mathcal{F}_s^{\mathcal{H}}(\boldsymbol{x}_s)\}), \{\boldsymbol{d}_i\}). \quad (3)$$

### 3.3. Enhancing Manifold Featurization

We have formulated the hash featurized manifold representation in Sec. 3.2 and this already delivers state-of-the-art NVS quality on challenging cross-scale scenarios. In this section, we further introduce two enhancements tailored for our representation, namely surface multisampling and manifold deformation, to better represent cross-scale details.

**Surface Multisampling.** Given the cross-scale, in-the-wild observations of a general large scene, casting a single ray per pixel neglects the unstructured scale variations and leads to blur or aliasing artifacts, due to the discrepancies in pixel colour when observing a surface point across varying distances or resolutions. To this end, we introduce a multisampling scheme for our surface-based featurization. Inspired by [5, 12, 38], we cast multiple rays per pixel to obtain a set of surface intersections $\{\boldsymbol{x}_s^{(j)}\}_{j=1}^{\gamma^2}$. To do so, we rasterize a $\gamma H \times \gamma W$ image, where each pixel in the original $H \times W$ image is super-sampled with a grid of $\gamma^2$ pixels. We then aggregate the information of the multiple surface intersections by individually querying the multi-resolution hash feature for each sample and pooling them with the mean operation:

$$\mathcal{F}_s^{\mathcal{H}}(\boldsymbol{x}_s) = \sum_{j=1}^{\gamma^2} \mathcal{F}_s^{\mathcal{H}}(\boldsymbol{x}_s^{(j)})/\gamma^2. \quad (4)$$

**Manifold Deformation.** Since our method directly featurizes a mesh, any geometric imperfections on the mesh will hinder the expressivity on local details. Similar to [39], we propose to further strengthen the multi-view consistency by a latent-space deformation $\xi : \mathbb{R}^{LZ} \times \mathbb{R}^3 \mapsto \mathbb{R}^{LZ}$. Specifically, we first featurize the surface using another hash encoding $\mathcal{F}_s^{\mathcal{D}}(\cdot)$ with learnable hash table $\mathcal{D}$. Then, a tiny MLP $\xi$ takes as inputs the new hash features and the view direction vector $\boldsymbol{d} \in \mathbb{R}^3$ to deform the initial surface in high-dimensional feature space:

$$\mathcal{F}_s^{\mathcal{H}'}(\boldsymbol{x}_s^{(j)}) = \mathcal{F}_s^{\mathcal{H}}(\boldsymbol{x}_s^{(j)}) + \xi\big(\mathcal{F}_s^{\mathcal{D}}(\boldsymbol{x}_s^{(j)}), \boldsymbol{d}\big). \quad (5)$$

As shown in Fig. 3 (b), equipped with surface multisampling and manifold deformation, our hash featurized manifold essentially represents a deformable frustum near the initial surface, making it more robust to handle scale variations and more flexible at capturing micro-scale details.

Figure 4. Illustration of the real-captured, unstructured, cross-scale imagery in our `GigaNVS` dataset. We collect high quality multi-view images at varying distances ranging from $5m$ to $10^3 m$.

## 4. Experiments

**Baselines.** We compare our method against explicit surface based representations [16, 39] and implicit volumetric representations [21, 25] with the same MVS mesh as additional overhead. For 3DGS [16], we initialize a dense set of gaussians on the mesh vertices and maintain as many gaussians as possible to make full use of the memory. For Zip-NeRF [5] and iNGP [25], we render the mesh into per-view depth map and supervise the volume-rendered depth similar to [7]. For Neuralangelo [21], we directly supervise the 3D zero-level set by the mesh vertices as done in [9]. The hash encoding parameters in [5, 21, 25] are set the same as ours. The implementation details can be found in the supplement.

### 4.1. `GigaNVS` Dataset

We introduce a novel dataset, namely `GigaNVS`, to evaluate our method and the baselines on the challenging task, i.e. cross-scale NVS of real-world large-scale scenes. This dataset contains 7 scenes of areas ranging from $1.3 \times 10^4 m^2$ to $3 \times 10^6 m^2$. For each scene, we collect thousands of high resolution (5K or 8K) multi-view images at multiple varying distances ($5m \sim 10^3 m$) using both aerial and ground photography. The camera poses are calculated by Agisoft Metashape [1] with subpixel-level reprojection errors. We refer readers to the supplement for more details.

An overview of our dataset is illustrated in Fig. 4 and a comparison to existing real-world multi-view datasets is listed in Table 1. To the best of our knowledge, `GigaNVS` is the first dataset characterized by gigantic scene scale and real-captured, cross-scale, high-resolution multi-view data. It fully exposes the scalability issues of state-of-the-art NVS algorithms, as verified in Sec. 4.2, and conveys unprecedentedly rich contents of real-world large scenes, which are rarely valued in existing datasets.

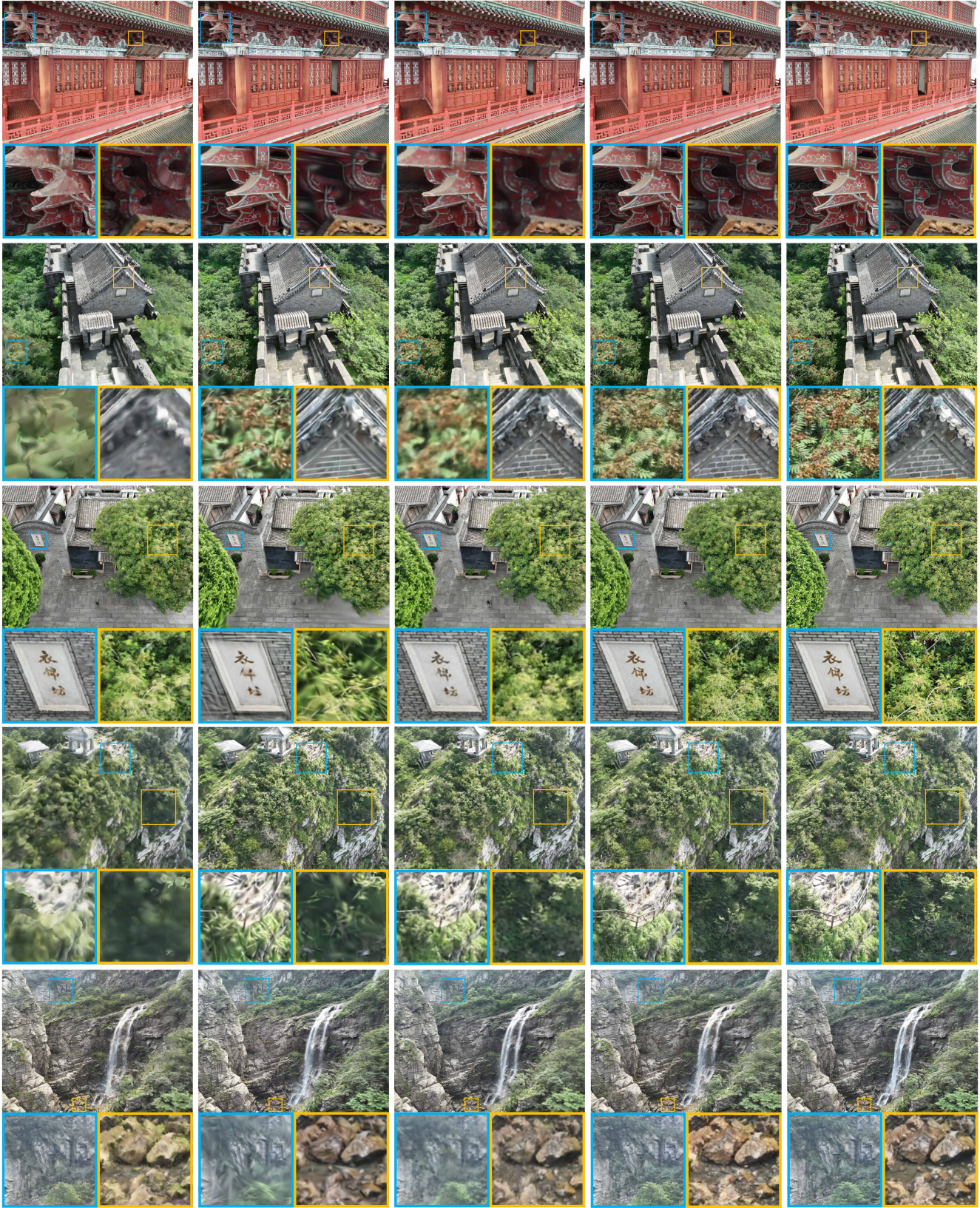| Datasets | #Scenes | #Images | Area ($m^2$) | Resol. | Cross-scale |
|---|---|---|---|---|---|
| Tanks&Temples [17] | 21 | $\sim 400$ | $\sim 10^2$ | 2K | $\times$ |
| ETH3D [33] | 25 | $\sim 40$ | $\sim 10^2$ | 6K | $\times$ |
| GigaMVS [51] | 13 | $\sim 300$ | $\sim 10^4$ | 19K | $\times$ |
| MipNeRF360 [4] | 7 | $\sim 200$ | $\sim 10^1$ | 3K/5K | $\times$ |
| **GigaNVS** (Ours) | 7 | $\sim 7000$ | $\sim 10^6$ | 5K/8K | $\checkmark$ |

Table 1. Comparison of real-world multi-view datasets, where #Images and Area denote the mean values across all scenes.

### 4.2. Comparative Results

**Benchmark on `GigaNVS` Dataset[†].** In Table 2 we report mean PSNR, SSIM [43], and LPIPS [52] metrics across the test views in `GigaNVS` dataset. Our method outperforms all state-of-the-art neural rendering methods by a large margin, and remarkably, achieves a 40% reduction in LPIPS relative to the second best method, ZipNeRF [5], demonstrating the significantly better perceptual fidelity of our method. The qualitative comparisons are shown in Fig. 5. Note that Meta representation [39] often produces blur and stretched artifacts due to the distortions of the UV-based featurization. 3DGS [16] delivers realistic rendering only at a global scale yet fails to capture the intricate details in close-ups. Implicit volumetric representations [5, 21, 25] suffer from surface ambiguities and struggle to handle complex large-scale structures even with geometric supervision, resulting in excessive blurries. By contrast, our method fully exploits the richness of the original imagery, enabling highly detailed NVS that is nearly indistinguishable from the ground truth.

**Benchmark on Tanks&Temples Dataset.** We compare against the state-of-the-arts on seven scenes from the Tanks&Temples dataset [17], and we report the mean metrics across all selected scenes in Table 3. As observed, our

---

[†]Since the baselines can not directly consume high resolution (5K/8K) inputs due to memory issues, we perform fair comparisons using downsampled images (1K) throughout the experiments. Please refer to the supplement for the high-resolution rendering of our method.

(a) Meta Representation [39]　　(b) 3D Gaussian Splatting [16]　　(c) ZipNeRF [5]　　(d) **Ours**　　(e) Ground Truth Image

Figure 5. Novel view synthesis results on the `GigaNVS` dataset. Compared to [5, 16, 39], our method robustly synthesizes realistic colour and intricate details, preserving approximately the input-level resolution. Please zoom-in to see the details.

| Scene | Meta Representation [39] | | | 3DGS [16] | | | Neuralangelo [21] | | | ZipNeRF [5] | | | Ours | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | PSNR ↑ | SSIM ↑ | LPIPS↓ | PSNR ↑ | SSIM ↑ | LPIPS↓ | PSNR ↑ | SSIM ↑ | LPIPS↓ | PSNR ↑ | SSIM ↑ | LPIPS↓ | PSNR ↑ | SSIM ↑ | LPIPS↓ |
| TW-Pavilion (Day) | 20.63 | 0.668 | 0.220 | 21.52 | 0.752 | 0.225 | 18.31 | 0.553 | 0.394 | 21.26 | 0.691 | 0.260 | 23.02 | 0.786 | 0.113 |
| TW-Pavilion (Night) | 23.82 | 0.750 | 0.235 | 24.38 | 0.816 | 0.199 | 21.73 | 0.671 | 0.374 | 24.59 | 0.800 | 0.198 | 25.55 | 0.845 | 0.120 |
| Lanes & Alleys | 19.39 | 0.734 | 0.188 | 20.30 | 0.787 | 0.222 | 18.10 | 0.633 | 0.321 | 20.39 | 0.796 | 0.169 | 20.41 | 0.814 | 0.110 |
| The Great Wall (T3) | 16.94 | 0.458 | 0.439 | 18.11 | 0.642 | 0.389 | 14.25 | 0.196 | 0.731 | 18.19 | 0.651 | 0.325 | 18.23 | 0.685 | 0.211 |
| The Great Wall (T2) | 19.86 | 0.701 | 0.226 | 19.91 | 0.709 | 0.314 | 19.01 | 0.635 | 0.319 | 18.99 | 0.780 | 0.202 | 21.03 | 0.806 | 0.131 |
| The Five Old Peaks | 18.98 | 0.614 | 0.330 | 18.79 | 0.750 | 0.272 | 16.65 | 0.417 | 0.555 | 19.93 | 0.741 | 0.251 | 20.05 | 0.774 | 0.146 |
| Sandie Spring | 16.64 | 0.546 | 0.409 | 17.25 | 0.670 | 0.380 | 13.98 | 0.264 | 0.678 | 17.32 | 0.679 | 0.298 | 17.61 | 0.724 | 0.206 |
| Mean | 19.47 | 0.639 | 0.292 | 20.04 | 0.732 | 0.286 | 17.43 | 0.481 | 0.482 | 20.09 | 0.734 | 0.243 | 20.84 | 0.776 | 0.148 |

Table 2. Quantitative comparisons on the GigaNVS dataset. Our method outperforms state-of-the-art approaches on all evaluation metrics.

method also demonstrates superiority on small-scale scenes from the public benchmark. Please refer to the supplement for more detailed metrics reporting and visual comparisons.

| Methods | PSNR↑ | SSIM↑ | LPIPS↓ |
|---|---|---|---|
| Meta Representation [39] | 28.32 | 0.892 | 0.119 |
| 3DGS [16] | 28.62 | 0.903 | 0.123 |
| Neuralangelo [21] | 27.41 | 0.898 | 0.114 |
| iNGP [25] | 28.67 | 0.905 | 0.110 |
| **Ours** | 29.66 | 0.914 | 0.079 |

Table 3. Quantitative evaluations on the Tanks&Temples dataset.

| Methods | PSNR↑ | SSIM↑ | LPIPS↓ |
|---|---|---|---|
| w/o Deformation | 22.77 | 0.775 | 0.131 |
| w/o Multisampling | 22.24 | 0.752 | 0.140 |
| Full Implementation | 23.02 | 0.786 | 0.113 |

Table 4. Ablation on featurization enhancements.

| Method | #Primitives | PSNR↑ | SSIM↑ | LPIPS↓ |
|---|---|---|---|---|
| Ours | 10M | 23.02 | 0.786 | 0.113 |
| Ours | 1M | 22.84 | 0.782 | 0.114 |
| 3DGS [16] | 3M | 21.52 | 0.752 | 0.225 |

Table 5. Ablation on mesh resolution.

| Methods | Time | | GPU Mem. |
|---|---|---|---|
| | Optimization | Rendering | |
| Meta Representation [39] | 35 h | 0.34 s | 20 GB |
| 3DGS [16] | 1 h | 0.01 s | 18 GB |
| ZipNeRF [5] | 8.3 h | 14.3 s | 19 GB |
| Neuralangelo [21] | 13 h | 32.2 s | 19 GB |
| iNGP [25] | 0.8 h | 1.19 s | 8 GB |
| Ours | 6.5 h | 0.08 s | 15 GB |

Table 6. Comparison of time and memory cost.

**Ablations on Featurization Enhancements.** We ablate the surface multisampling and manifold deformation module in our pipeline using the TW-Pavilion (Day) scene. As shown in Table 4, both designs improve the rendering quality.

**Ablations on Mesh Resolution.** In Table 5, we demonstrate the robustness of our method w.r.t. the mesh resolution, where mesh decimation [11] is used to obtain the mesh with desired face counts. Notably, our method achieves less than 1% of differences in all metrics when down-samping the mesh by 90%, indicating the effectiveness of our featurization without heavy reliance on the 3D discretization resolution. Our representation with 1 million triangle faces significantly outperforms the state-of-the-art 3DGS [16] with 3 million gaussians. Metrics are reported on the TW-Pavilion (Day) scene, and we refer readers to the supplement for more ablations on other scenes.

**Efficiency.** We compare the time and memory cost of the competing methods in Table 6, using the TW-Pavilion (Day) scene. The rendering time and GPU memory cost are measured by synthesizing an image of 989×1320 resolution. Our method is significantly faster than volume rendering based methods [5, 21, 25] due to the rasterization pipeline. Compared to [39], our method also shows superior efficiency by using a light-weight neural shader. Note that 3DGS [16] demonstrates the highest time efficiency leveraging the splatting and SH-based calculations, whereas iNGP [25] and ours are more efficient in memory by incorporating compact neural components.

## 5. Conclusion

We introduce hash featurized manifold, a novel representation for high-fidelity cross-scale neural rendering of real-world large-scale scenes. The core is an expressive surface-based featurization constructed by guiding the volumetric hash encoding with the rasterization of a surface manifold. Our representation fully unlocks the expressivity of multi-resolution hash encoding by skipping the redundant space and concentrating on multi-view consistent colour gradient. We also propose a novel GigaNVS dataset consisting of seven real-world large scenes to benchmark cross-scale, high-resolution novel view synthesis. Extensive experiments demonstrate that our method achieves unparalleled levels of realism by effectively reflecting both macro-scale and micro-scale scene contents.

**Limitation & Future Work** Although showing robustness to the mesh resolution, our method currently can not handle the incompleteness and occlusions caused by the incorrect geometry. Our future work is to exploit differentiable rendering for more flexible control of the geometry.

# References

[1] Agisoft LLC. Agisoft metashape, 2021. 6

[2] Kara-Ali Aliev, Artem Sevastopolsky, Maria Kolos, Dmitry Ulyanov, and Victor Lempitsky. Neural point-based graphics. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXII 16*, pages 696–712. Springer, 2020. 2, 3, 4

[3] 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

[4] 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. 1, 6

[5] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Zip-nerf: Anti-aliased grid-based neural radiance fields. *arXiv preprint arXiv:2304.06706*, 2023. 1, 2, 3, 4, 5, 6, 7, 8

[6] 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. 1

[7] 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. 6

[8] Michael S Floater and Kai Hormann. Surface parameterization: a tutorial and survey. *Advances in multiresolution for geometric modelling*, pages 157–186, 2005. 2, 4

[9] Qiancheng Fu, Qingshan Xu, Yew Soon Ong, and Wenbing Tao. Geo-neus: Geometry-consistent neural implicit surfaces learning for multi-view reconstruction. *Advances in Neural Information Processing Systems*, 35:3403–3416, 2022. 6

[10] Yasutaka Furukawa, Carlos Hernández, et al. Multi-view stereo: A tutorial. *Foundations and Trends® in Computer Graphics and Vision*, 9(1-2):1–148, 2015. 2

[11] Michael Garland and Paul S Heckbert. Surface simplification using quadric error metrics. In *Proceedings of the 24th annual conference on Computer graphics and interactive techniques*, pages 209–216, 1997. 8

[12] Ned Greene and Paul S Heckbert. Creating raster omnimax images from multiple perspective views using the elliptical weighted average filter. *IEEE Computer Graphics and Applications*, 6(6):21–27, 1986. 5

[13] Jiaming Gu, Minchao Jiang, Hongsheng Li, Xiaoyuan Lu, Guangming Zhu, Syed Afaq Ali Shah, Liang Zhang, and Mohammed Bennamoun. UE4-neRF:neural radiance field for real-time rendering of large-scale scene. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. 1, 3

[14] Kai Hormann, Bruno Lévy, and Alla Sheffer. Mesh parameterization: theory and practice. *ACM SIGGRAPH ASIA 2008 courses*, 2007. 2, 4

[15] Mengqi Ji, Juergen Gall, Haitian Zheng, Yebin Liu, and Lu Fang. Surfacenet: An end-to-end 3d neural network for multiview stereopsis. In *Proceedings of the IEEE international conference on computer vision*, pages 2307–2315, 2017. 2

[16] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics (ToG)*, 42(4):1–14, 2023. 1, 2, 3, 4, 6, 7, 8

[17] 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. 1, 2, 6

[18] Georgios Kopanas, Thomas Leimkühler, Gilles Rainer, Clément Jambon, and George Drettakis. Neural point catacaustics for novel-view synthesis of reflections. *ACM Transactions on Graphics (TOG)*, 41(6):1–15, 2022. 2, 3, 4

[19] Yixuan Li, Lihan Jiang, Linning Xu, Yuanbo Xiangli, Zhenzhi Wang, Dahua Lin, and Bo Dai. Matrixcity: A large-scale city dataset for city-scale neural rendering and beyond. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3205–3215, 2023. 3

[20] Zhuopeng Li, Lu Li, and Jianke Zhu. Read: Large-scale neural scene rendering for autonomous driving. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 1522–1529, 2023. 3

[21] Zhaoshuo Li, Thomas Müller, Alex Evans, Russell H Taylor, Mathias Unberath, Ming-Yu Liu, and Chen-Hsuan Lin. Neuralangelo: High-fidelity neural surface reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8456–8465, 2023. 2, 3, 4, 6, 8

[22] Liqiang Lin, Yilin Liu, Yue Hu, Xingguang Yan, Ke Xie, and Hui Huang. Capturing, reconstructing, and simulating: the urbanscene3d dataset. In *ECCV*, pages 93–109, 2022. 3

[23] Jeffrey Yunfan Liu, Yun Chen, Ze Yang, Jingkang Wang, Sivabalan Manivasagam, and Raquel Urtasun. Real-time neural rasterization for large scenes. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8416–8427, 2023. 1, 2, 3, 4

[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. *Communications of the ACM*, 65(1):99–106, 2021. 2, 3

[25] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)*, 41(4):1–15, 2022. 2, 3, 4, 6, 8

[26] Alexander Raistrick, Lahav Lipson, Zeyu Ma, Lingjie Mei, Mingzhe Wang, Yiming Zuo, Karhan Kayan, Hongyu Wen, Beining Han, Yihan Wang, et al. Infinite photorealistic worlds using procedural generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12630–12641, 2023. 3

[27] Ruslan Rakhimov, Andrei-Timotei Ardelean, Victor Lempitsky, and Evgeny Burnaev. Npbg++: Accelerating neural point-based graphics. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15969–15979, 2022. 2, 3, 4

[28] Nicolas Ray, Wan-Chiu Li, Bruno Lévy, Alla Sheffer, and Pierre Alliez. Periodic global parameterization. *ACM Trans. Graph.*, 25:1460–1485, 2006. 2, 4

[29] Konstantinos Rematas, Andrew Liu, Pratul P Srinivasan, Jonathan T Barron, Andrea Tagliasacchi, Thomas Funkhouser, and Vittorio Ferrari. Urban radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12932–12942, 2022. 1, 2, 3

[30] Radu Alexandru Rosu and Sven Behnke. Permutosdf: Fast multi-view reconstruction with implicit surfaces using permutohedral lattices. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8466–8475, 2023. 2, 3, 4

[31] Darius Rückert, Linus Franke, and Marc Stamminger. Adop: Approximate differentiable one-pixel point rendering. *ACM Transactions on Graphics (ToG)*, 41(4):1–14, 2022. 2, 3, 4

[32] Johannes L Schönberger, Enliang Zheng, Jan-Michael Frahm, and Marc Pollefeys. Pixelwise view selection for unstructured multi-view stereo. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part III 14*, pages 501–518. Springer, 2016. 2

[33] Thomas Schops, Johannes L Schonberger, Silvano Galliani, Torsten Sattler, Konrad Schindler, Marc Pollefeys, and Andreas Geiger. A multi-view stereo benchmark with high-resolution images and multi-camera videos. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3260–3269, 2017. 1, 6

[34] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5459–5469, 2022. 3

[35] Matthew Tancik, Vincent Casser, Xinchen 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. 1, 2, 3

[36] Justus Thies, Michael Zollhöfer, and Matthias Nießner. Deferred neural rendering: Image synthesis using neural textures. *Acm Transactions on Graphics (TOG)*, 38(4):1–12, 2019. 2, 3, 4

[37] 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. 1, 2, 3

[38] Chen Wang, Xian Wu, Yuan-Chen Guo, Song-Hai Zhang, Yu-Wing Tai, and Shi-Min Hu. Nerf-sr: High quality neural radiance fields using supersampling. In *Proceedings of the*

[39] G. Wang, J. Zhang, K. Zhang, R. Huang, and L. Fang. Giganticnvs: Gigapixel large-scale neural rendering with implicit meta-deformed manifold. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, (01):1–15, 2023. 1, 2, 3, 4, 5, 6, 7, 8

[40] 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. *arXiv preprint arXiv:2106.10689*, 2021. 2, 3

[41] Xueyang Wang, Xiya Zhang, Yinheng Zhu, Yuchen Guo, Xiaoyun Yuan, Liuyu Xiang, Zerun Wang, Guiguang Ding, David Brady, Qionghai Dai, et al. Panda: A gigapixel-level human-centric video dataset. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3268–3278, 2020. 3

[42] Yiming Wang, Qin Han, Marc Habermann, Kostas Daniilidis, Christian Theobalt, and Lingjie Liu. Neus2: Fast learning of neural implicit surfaces for multi-view reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3295–3306, 2023. 2, 3, 4

[43] 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. 6

[44] 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 *European conference on computer vision*, pages 106–122. Springer, 2022. 2, 3

[45] Bangbang Yang, Chong Bao, Junyi Zeng, Hujun Bao, Yinda Zhang, Zhaopeng Cui, and Guofeng Zhang. Neumesh: Learning disentangled neural mesh-based implicit field for geometry and texture editing. In *European Conference on Computer Vision*, pages 597–614. Springer, 2022. 2, 3, 4

[46] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. In *Proceedings of the European conference on computer vision (ECCV)*, pages 767–783, 2018. 2

[47] 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. 3

[48] 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. 3

[49] Xiaoyun Yuan, Mengqi Ji, Jiamin Wu, David J Brady, Qionghai Dai, and Lu Fang. A modular hierarchical array camera. *Light: Science & Applications*, 10(1):1–9, 2021. 3

[50] Jianing Zhang, Tianyi Zhu, Anke Zhang, Xiaoyun Yuan, Zihan Wang, Sebastian Beetschen, Lan Xu, Xing Lin, Qionghai Dai, and Lu Fang. Multiscale-vr: multiscale gigapixel 3d panoramic videography for virtual reality. In *2020 IEEE international conference on computational photography (ICCP)*, pages 1–12. IEEE, 2020. 3

30th ACM International Conference on Multimedia, pages 6445–6454, 2022. 5

[51] Jianing Zhang, Jinzhi Zhang, Shi Mao, Mengqi Ji, Guangyu Wang, Zequn Chen, Tian Zhang, Xiaoyun Yuan, Qionghai Dai, and Lu Fang. Gigamvs: a benchmark for ultra-large-scale gigapixel-level 3d reconstruction. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(11):7534–7550, 2021. 1, 3, 6

[52] 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. 6

[53] Yiming Zuo and Jia Deng. View synthesis with sculpted neural points. *arXiv preprint arXiv:2205.05869*, 2022. 2, 3, 4