

# QFF: Quantized Fourier Features for Neural Field Representations

Jae Yong Lee<sup>1</sup>, Yuqun Wu<sup>1</sup>, Chuhan Zou<sup>\*2</sup>, Shenlong Wang<sup>1</sup>, and Derek Hoiem<sup>1</sup>

<sup>1</sup>University of Illinois at Urbana-Champaign

<sup>2</sup>Amazon.com

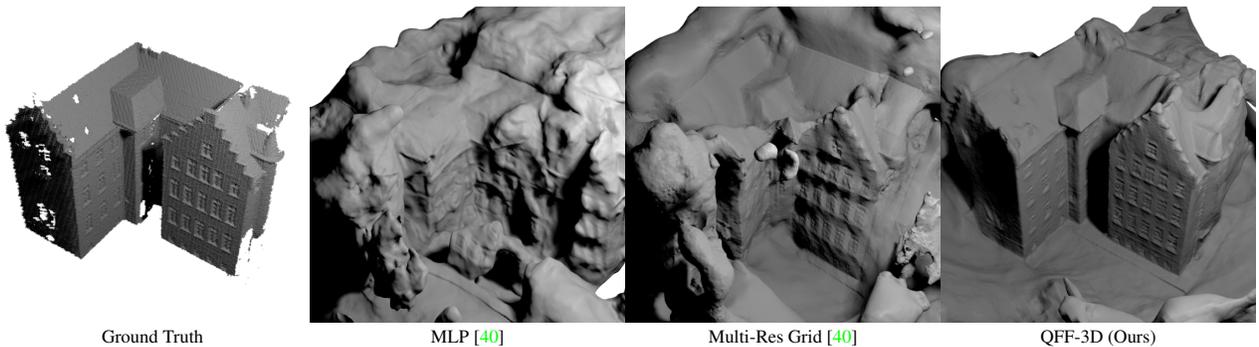


Figure 1. We propose Quantized Fourier Features (QFF), an easy-to-train, memory-efficient yet expressive neural field representation. We compare the multilayer perceptron (MLP), multi-resolution grids, and our proposed QFF for neural surface reconstruction over sparse, 3-view reconstructions in the DTU Dataset [1]. The QFF model encodes high-frequency details better than the baseline MLP and preserves smoothness where appropriate, resulting in better geometry compared to the multi-resolution grids.

## Abstract

*Multilayer perceptrons (MLPs) learn high frequencies slowly. Recent approaches encode features in spatial bins to improve speed of learning details, but at the cost of larger model size and loss of continuity. Instead, we propose to encode features in bins of Fourier features that are commonly used for positional encoding. We call these Quantized Fourier Features (QFF). As a naturally multiresolution and periodic representation, our experiments show that using QFF can result in smaller model size, faster training, and better quality outputs for several applications, including Neural Image Representations (NIR), Neural Radiance Field (NeRF) and Signed Distance Function (SDF) modeling. QFF are easy to code, fast to compute, and serve as a simple drop-in addition to many neural field representations.*

## 1. Introduction

Deep networks’ impressive ability as function approximators has been demonstrated, e.g. for physics [24], compression [29], and 3D modeling [16, 19]. Yet, learning

high-frequency functions is still a challenge. Deep networks tend to learn low-frequency signals first, and higher frequencies later, requiring long training times. This phenomenon, known as spectral bias [23] or frequency principle (*f-principle*) [36], is universally observed in multilayer perceptron (MLP) architectures. For instance, optimizing Neural Radiance Field (NeRF) [16] MLP models takes hours, or even days for large scenes [16]. The success of NeRF relies on encoding position with multi-frequency sinusoidal coefficients [16, 31] that are amenable to linear functions, but, while quality greatly improves, convergence is still slow because changes to individual parameters can have widespread effects. Much faster convergence can be achieved by partitioning the function’s domain into local components, e.g. voxels in a 3D scene, with parameters that can be independently optimized [17, 25]. However, spatial gridding reduces continuity and requires multi-resolution grids with many parameters.

Our **key insight** is that quantizing the multi-frequency sinusoidal coefficients provides a naturally multiresolution representation that enables the high quality of MLP models and fast training of spatial quantization approaches, with minimal computational and memory cost. A feature vector is learned for each quantized value of a coefficient, and these features are added to the continuous coefficients as

\*The work is not related to the author’s position at Amazon.

the input to an MLP. We call these quantized Fourier features (QFF). By quantizing multiple frequencies, the QFF maintain continuity and high spatial resolution without redundancy. When modeling functions with sparse outputs, such as occupancy in 3D scenes, the periodicity of QFF provides further advantage of encoding the full domain at high resolution without using many parameters to encode completely empty space. QFF can be computed per input dimension (*QFF-Lite*), or factorized and composed into multiple 2D and 1D grids (*QFF-3D*) to encode more detail and further speed convergence. QFF is simple to code, fast to compute, and easily inserted into many neural field frameworks. Our experiments demonstrate advantages in convergence and quality of models in applications of image encoding, novel view synthesis, and 3D surface modeling.

In summary, our **contribution** is a quantized Fourier feature (QFF) that is easy to compute and include in neural representations and provides advantages of small model size, high resolution, fast training, and excellent model quality across several applications.

## 2. Related Work

Neural networks tend to learn low frequencies of the objective function before high frequencies. This phenomenon is known as spectral bias [23] or the Frequency Principle [36], and has been shown to hold for synthetic and real data, as well as multilayer perceptrons, convolutional networks, and other architectures [37]. This has direct implications for learning Neural Field Representations (NFR), manifested in a variety of applications including physics-inspired neural networks [3, 24], shape fitting [4, 19], image compression [27], and neural radiance fields [2, 15, 16, 28]. NFR require longer convergence time when learning high-dimensional, high-frequency functions, which is the main bottleneck for large-scale training.

Multiple approaches has been proposed to embed feature vectors into explicit geometric structures (e.g. voxels) [13, 14, 17, 25, 39]. This strategy typically achieve faster convergence (compared to MLP), they inherently use more memory [17] and are known to overfit [34, 35]. A common strategy to avoid overfitting is to employ gradient-based regularization [8, 9] to simulate smoothness. To reduce memory for storing explicit mappings, sparse-grid [13] and octree [30] are proposed to allocate more parameters near occupied scene regions. Spatial Hashing [17] uses a multidimensional hash table to map input coordinates to random hash table indices. This saves storage but creates a discontinuous representation, which leads to artifacts in inherently smooth functions such as signed distance fields (SDF) [40]. Approaches such as TensorRF [6] reduce parameters by decomposing multi-dimensional arrays into compositions of lower-dimensional arrays. Many of these approaches require knowledge of a bounding volume or occupied portions

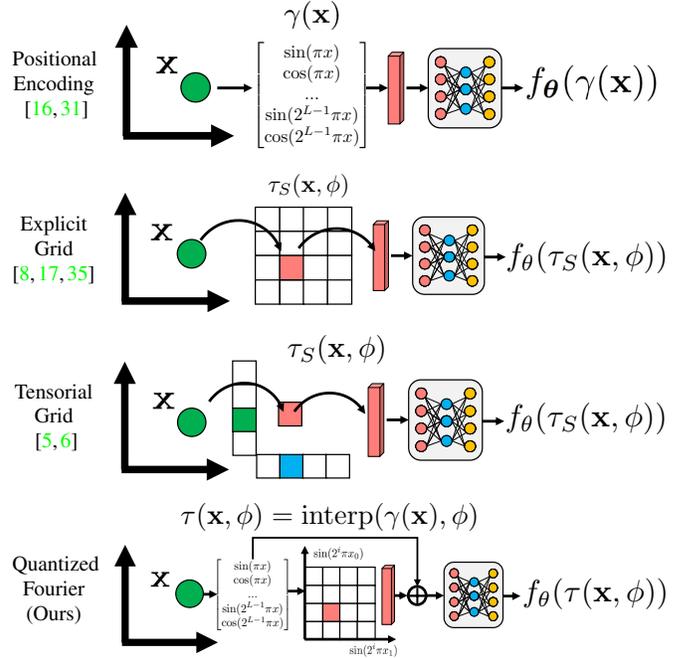


Figure 2. **Comparison between different neural field representations.** From top to bottom: positional encoding, explicit spatial grids, tensorial grids, and our proposed quantized Fourier features.

of the scene, necessitating a coarse-to-fine training strategy is used to first identifying volume of interest in coarse volume, and crop-and-resize to fit the finer scenes.

Fig. 2 compares several representations. Our QFF representation uses explicit geometric structure, but in the form of quantized Fourier components so that it is continuous, compact, and multiscale, and does not require prior knowledge of the bounding box or occupied space. We represent 3D volumes, either factorized per dimension (*QFF-Lite*) or by composing 2D arrays (*QFF-3D*), similar to TensorRF. Further, the QFF representation is based on commonly used positional encodings and easily inserted into most Neural Field Representations.

## 3. Background: Neural Field Representations

### 3.1. Neural Fields with Positional Encodings

Neural field representations approximate the mapping from coordinates  $\mathbf{x} \in \mathbb{R}^K$  to signal values  $\mathbf{v} \in \mathbb{R}^D$  (e.g., color or opacity) with a learnable neural network  $f_\theta(\cdot)$ , often instantiated with a multilayer perceptron (MLP). An essential recipe for neural fields is positional encoding [16, 31], which maps the Euclidean coordinates input to sinusoidal activations across  $L$  different frequency levels:

$$\gamma(x_k) = [\sin(2^{0 \dots L-1} \pi x_k), \cos(2^{0 \dots L-1} \pi x_k)] \quad (1)$$

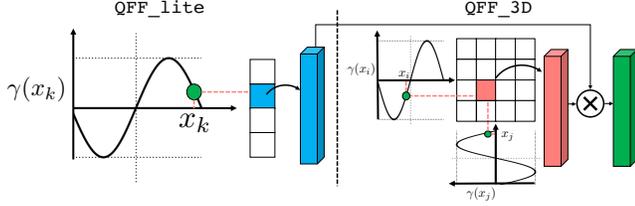


Figure 3. Visualization of single-dimensional QFF-lite and QFF-3D feature sampling procedure (applied separately for each positional encoding component). QFF-lite: Each positional encoding  $\gamma(x_k)$  is discretized into  $M$  bins. Each bin contains a multi-dimensional learnable feature  $\phi$ . The encoded value  $\gamma(x_k)$  is used as an index to query features from  $\phi$ :  $\tau(\gamma(x_k), \phi)$ . QFF-3D: We query QFF-lite features  $\tau(\gamma(x_k), \phi)$ . We then build  $M \times M$  bins and use the positional encoding of  $(\gamma(x_i), \gamma(x_j))$  to query the corresponding features using bilinear interpolation. We compute the element-wise product between  $\tau(x_i, x_j, \phi)$  and  $\tau(\gamma(x_k), \phi)$  to get the final feature.

$\gamma(\mathbf{x}) = \text{concat}([\gamma(x_0) \dots \gamma(x_{K-1})])$  is a continuous, multi-scale, periodic representation of  $\mathbf{x}$  along each coordinate. Then,  $\mathbf{v}$  is computed as:

$$\mathbf{v} = f_{\theta}(\gamma(\mathbf{x})) \quad (2)$$

Neural fields representations with positional encodings can capture high-frequency structures in the scene, making them suitable for representing complex spatial-temporal signals across various modalities, such as images [27], videos [22], and 3D shapes [16, 33]. However, it takes a long time to converge, because all MLP parameters have to update for every pair of inputs and outputs.

### 3.2. Neural Fields with Explicit Features

An alternative paradigm [6, 8, 17, 21, 35, 41] is to explicitly encode features at the input position  $\tau_s(\mathbf{x}; \phi)$ , where  $\phi$  are learnable feature tensors that are linearly interpolated based on  $\mathbf{x}$  to get feature values. Often, the feature values are added to the positional encoding values as input to the network, which produces  $\mathbf{v}$ . Popular choices for  $\tau_s(\mathbf{x}; \phi)$  include trilinearly interpolated voxel features [35]; tri-plane representation [6, 21] and spatial voxel hashing [17]. Compared to the positional encoding-based neural field representation, such representations are significantly faster to train and evaluate. Furthermore, the neural fields are readily compositional and scalable. But because existing explicit feature representations are in the *spatial* domain, they are sensitive to resolution and may waste memory representing features in empty portions of space.

## 4. Quantized Fourier Features (QFF)

We present Quantized Fourier Features (QFF), an easy-to-train, memory-efficient yet expressive neural field representation. Our proposed approach combines the best worlds

```

import torch
import torch.nn.functional as F

# L: number of frequency levels
# K: input dimension
# M: quantization levels
# N: feature dimension
# Phi: frequency feature grid (Kx2L, M, N)
# x: input

def NeuralField(x, L):
    gamma = PositionalEncoding(x, L)
    y = MLP(gamma)
    return y

def QFF_lite(x, L, Phi):
    gamma = PositionalEncoding(x, L)
    tau = F.grid_sample(Phi, gamma)
    # Quantized Fourier Feature
    y = MLP(gamma + tau)
    return y

```

Figure 4. **Apply QFF to existing methods.** We show PyTorch [20] pseudo-code to apply QFF to existing MLP based systems, which can be implemented in less than 4 lines of code.

of explicit feature representation and frequency-based positional encoding. The key idea is to store explicit features in the Fourier domain and use quantized positional encoding to query the feature.

### 4.1. QFF-Light

We first explain QFF when encoding each input dimension separately. As shown in Figure 3, we create  $M$  bins for each  $i$ -th position encoding component from Eq. 1, and each bin stores a learnable  $N$ -dimensional feature vector. This gives us a  $K \times 2L \times M \times N$ -dimensional tensor  $\phi$ , where  $K$  is the input dimension (e.g., 3 for a volume),  $L$  is the number of frequencies,  $M$  is the quantization level, and  $N$  is feature dimension. Given the spatial input coordinate value  $x_k$ , we compute features for each  $i$ -th positional encoding value by linear interpolation:

$$\tau(x_k; \phi)_i = \text{interp}(\gamma(x_k)_i, \phi_{ki}) \quad (3)$$

where  $\gamma(x_k)_i$  denotes the  $i$ -th positional encoding value of  $x_k$  and  $\phi_{ki}$  is a  $M \times N$  slice of  $\phi$  after indexing with  $k$  and  $i$ . Figure 3 depicts the feature computation process.

As QFF and Multi-dimensional QFF are piecewise-linear functions, the derivative at the border of quantization in QFF changes sharply due to the discretization. Hence, we add the original positional encoding values to the quantized features to induce smoothness. We name the per-dimension QFF, which adds the positional encoding, as QFF-Lite:

$$\text{QFF-Lite}(x_k; \phi)_i = \tau(x_k, \phi)_i + \gamma(x_k)_i. \quad (4)$$

The resulting values are concatenated across  $k$  and  $i$ , input to the MLP, and trained via gradient descent.

## 4.2. QFF-3D

So far, QFF is computed independently for each coordinate. When applying the 1-D QFF embeddings to multi-dimensional input ( $n > 1$ ), the MLPs must learn to model correlations across multi-dimensional features. This brings an additional burden to MLPs and makes them less capable or slower to converge. We can instead store features along a grid or volume. For example, for 2-D QFF, we create  $M \times M$  bins for each positional encoding component across both dimensions, so that  $\phi$  is  $K \times 2L \times M^2 \times N$ -dimensional, and bilinearly interpolate at query time.

For 3-dimensional data, we use the multi-dimensional QFF by TensorRF [6] style of spatial decomposition. In particular, we follow vector-matrix (VM) decomposition and define QFF-3D as:

$$\text{QFF-3D}(x_0; x_1, x_2, \phi)_i = \tau(x_0, \phi)_i \cdot \tau(x_1, x_2; \phi)_i + \gamma(x_0)_i \quad (5)$$

Equation 5 describes the VM decomposition in  $x_0$  dimension. We apply the VM decomposition in all 3 dimensions, as done in TensorRF [6].

## 4.3. Application of QFF

QFF is fast to compute and differentiable. Unlike spatial gridding approaches, where 3D positions and features directly correspond, each spatial position in QFF is represented in multiple frequencies, and each frequency component is represented at multiple positions, as depicted in Figure 3. This enables the network to use the parameters to model occupied portions of the scene, without knowing which portions are occupied in advance, improving convergence and resolution for a given number of parameters.

QFF is also easy to plug and play into the existing neural fields that use Fourier positional encodings. The addition of QFF is independent of the flow of existing methods, and the computation of positional encodings can be reused when computing the QFF feature. Figure 4 illustrates how to apply QFF to existing methods in few lines of code to speed up training and improve performance.

## 5. Experiments

We experiment with our proposed QFF on a wide range of neural field representations:

- (1) **Neural Image Representations** (Section 5.1);
- (2) **Neural Radiance Fields** (Section 5.2);
- (3) **Neural Signed Distance Fields**, including few-shot object reconstruction and large-scale scene reconstruction (Section 5.3).

Method	Pos. Enc.	QFF-Lite	Natural	Text
ReLU	×	×	17.74	18.38
	✓	×	29.41	30.60
	✓	✓	<b>30.30</b>	<b>32.05</b>
SIREN [27]	×	×	28.18	30.94
	✓	×	37.06	<b>47.34</b>
	✓	✓	<b>37.68</b>	46.92

Table 1. **Evaluation on supervised Neural Image Representation in PSNR.** The best performing method for each dataset (Natural, Text) for each activation function (ReLU, SIREN) is marked in bold. Overall, we see improvements by applying *QFF-Lite*.

Through experiments on (1) and (2), we show our improvement over standard positional encoding and efficiency in spatial representation due to spatial hashing. We further demonstrate our ability to learn continuous representations in (3). We apply QFF-3D and QFF-Lite to existing state-of-the-art and baseline architectures, to demonstrate the improvement purely due to adding QFF to the representation.

### 5.1. Neural Image Representations

**Dataset.** We use natural and text mega-pixel images from Tancik et al. [31]. Natural images contain high frequency signals and text images contain sparse signals.

**Evaluation Metrics.** We evaluate the reconstruction ability by measuring the peak signal-to-noise ratio (PSNR).

**Implementation.** For 2D image representation, we use QFF-Lite with  $N = 1$ : a single value is added to each positional encoding value after sampling. We use  $L = 128$  for the number of scales and  $M = 2^7$  for the feature bin resolution. We use Adam [11] optimizer with  $lr = 5e^{-4}$  for both QFF and MLP parameters.

**Results.** We compare application of our QFF-Lite on ReLU and SIREN [27] activation functions, with and without our method for supervised 2D image reconstruction task. Table 1 shows the results for natural images and text images reconstruction. For both ReLU and SIREN, after applying QFF-Lite, we achieve a slight improvement (about 0.6) in PSNR on natural images. ReLU using QFF-Lite in text images achieves a significant improvement (1.45). We see a slight drop in PSNR for text images on SIREN, but we note that the PSNR of SIREN text reconstruction with positional encoding has peaked; there is not much room for improvement.

### 5.2. Neural Radiance Fields

**Dataset.** We use the NeRF Synthetic dataset [16] to compare with existing methods.

**Evaluation Metrics.** We evaluate the novel view synthesis performance by measuring the peak signal-to-noise ratio (PSNR).

Method		Chair	Drums	Ficus	Hotdog	Lego	Materials	Mic	Ship	Mean	Params	Steps
Decomp.	NeRF [16]	33.00	25.01	30.13	36.18	32.54	29.62	32.91	28.65	31.00	1.19 M	100+K
	TensorRF-CP [6]	<b>33.60</b>	25.17	30.72	<b>36.24</b>	34.05	30.10	33.77	<b>28.84</b>	31.56	728 K	<b>30 K</b>
	NeRFacc [12]	33.32	<b>25.39</b>	<b>32.52</b>	35.80	33.69	29.73	33.76	28.18	31.55	618 K	50 K
	NeRFacc (QFF-Lite)	33.36	25.33	31.97	35.70	<b>34.16</b>	<b>30.15</b>	<b>33.90</b>	28.20	<b>31.59</b>	<b>522 K</b>	50 K
Comp.	Instant-NGP [17]	35.00	<b>26.02</b>	33.51	37.40	36.39	29.78	36.22	<b>31.10</b>	33.18	12.6 M	50 K
	TensorRF-VM [6]	<b>35.76</b>	26.01	<b>33.99</b>	<b>37.41</b>	36.46	<b>30.12</b>	34.61	30.77	33.14	17.6 M	<b>30 K</b>
	NeRFacc (QFF-3D)	35.69	25.97	33.70	37.07	<b>36.68</b>	30.07	<b>37.53</b>	29.97	<b>33.35</b>	<b>9.82 M</b>	50 K

Table 2. **PSNR evaluation on test images of the NeRF Synthetic dataset.** The top and bottom rows show methods without and with feature composition, respectively. For the baseline methods, we use the numbers reported in original papers. We mark the best method for each section in bold. Our method achieves better mean image quality for novel view synthesis with fewer model parameters.

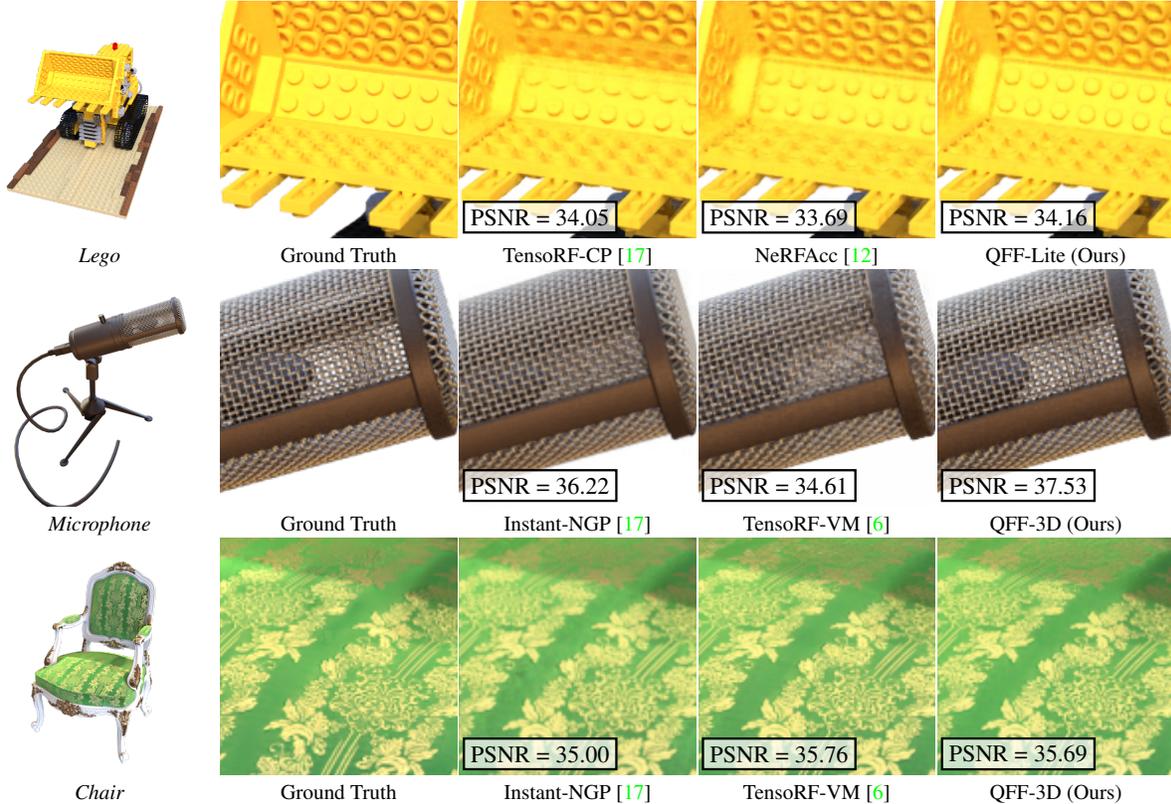


Figure 5. **Comparison of different neural radiation field representations.** The first row compares the decomposed representation of the scene *Lego*. The bottom two rows compare the composed representations of the scene *Microphone* and *Chair*. We obtain images of TensorRF [6] from the results provided by the authors, and NeRFacc [12] and Instant-NGP [17] results by running the code provided by the authors.

**Implementation.** We apply our QFF-Lite and QFF-3D to the Neural Radiance Field baseline NeRFacc [12]. Table 2 summarizes our results. We divide our results into decomposed scene representations, which do not explicitly compose input encodings (e.g. QFF-Lite), and composed scene representations, which explicitly create a composed 3D representation (e.g. QFF-3D). For both of our methods, we use  $N = 16$ ,  $L = 6$  and  $M = 2^7$ , and use Adam [11] optimizer with  $lr = 1e-2$  for the QFF and  $lr = 5e-4$  for

MLP. \* We use a 4-layer MLP for QFF-Lite and a 2-layer MLP for QFF-3D for all scenes.

**Results.** As shown in Table 2, in methods using decomposed scene representations, our QFF-Lite achieves the highest overall PSNR with smaller number of parameters,

\*For scene *Chair* and *Ficus*, we use slightly lower learning rate of  $9.5e-3$  and  $3.5e-4$  for the QFF and the MLP due to the loss converging to NaN in NeRFacc [12].

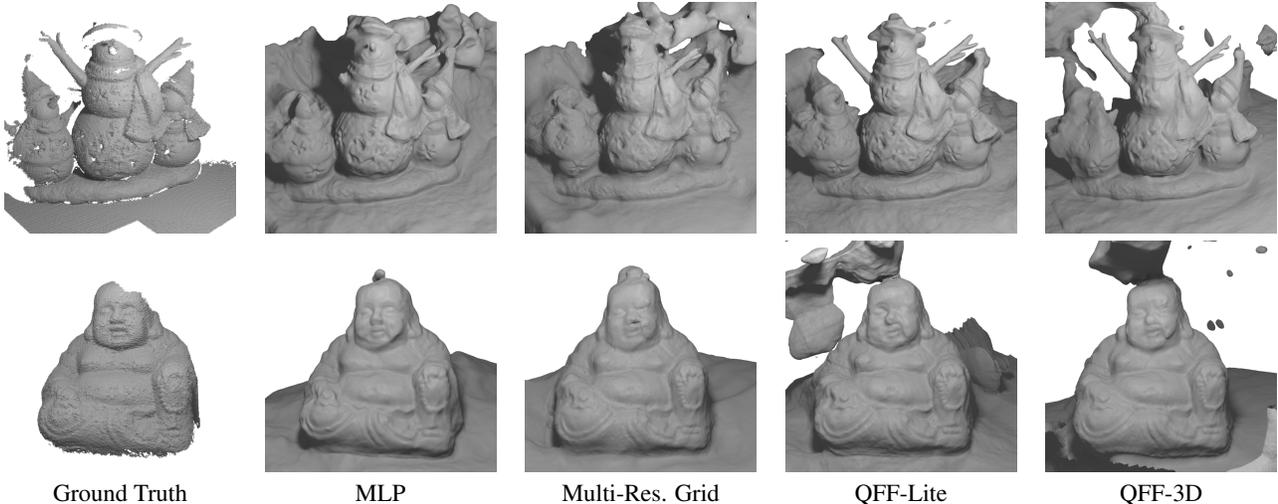


Figure 6. **Qualitative visualization of DTU dataset [1] reconstruction based on 3 sparse views.** MLP and Multi-Res Grid refer to the MonoSDF [40] with corresponding architectures. We use author provided meshes for visualizing results with MLP and Multi-Res Grid.

Method	Chamfer- $L_1$	# Params
TSDF-Fusion [26]	4.80	-
COLMAP [26]	2.56	-
RealityCapture	2.84	-
MonoSDF (Multi-Res. Grids) [40]	3.68	12.5 M
MonoSDF (MLP) [40]	<b>1.86</b>	<b>670 K</b>
MonoSDF (QFF-Lite)	2.15	837 K
MonoSDF (QFF-3D)	2.79	5.10 M

Table 3. **Evaluation on DTU dataset [1] with 3 views.** We report the Chamfer- $L_1$  distance of baseline and our results. We mark the best performing method in bold.

compared to all baseline methods. For methods with composed scene representations, our QFF-3D results are mostly on par with existing methods, with PSNR differences of at most 0.1 for 3 out of 8 scenes, and achieve slightly higher average PSNR with a smaller number of parameters. We emphasize that we do not impose a total variation (TV) loss or have to crop and resize the feature vectors as in TensorRF [6]. This is because our representation adds original positional encoding to preserve smoothness and is naturally multiscale and shift-invariant.

We present our qualitative comparisons in Figure 5. In both composed and decomposed scene representations, our model is able to capture high frequency details, such as inner-microphone geometry and Lego textures, with larger PSNR improvements compared to baseline models.

### 5.3. Neural Signed Distance Fields

**Datasets.** We use DTU dataset [1] for sparse view (3-view) object reconstruction and ScanNet dataset [7] for

large-scale reconstruction.

**Evaluation Metrics.** For the DTU dataset [1], we measure the Chamfer distance using the evaluation protocol provided by the dataset. For the ScanNet dataset, we report the Chamfer distance and  $F_1$  score following the baseline methods [10, 40].

**Sparse view object reconstruction** We apply our QFF-Lite and QFF-3D on MonoSDF [40] for DTU dataset [1] with three views. We use  $N = 8$ ,  $M = 2^7$   $L = 6$ , 8-layer MLP for QFF-Lite and 2-layer MLP for QFF-3D for all scenes. Table 3 shows the results of our methods compared to the baselines. We find that both our QFF-Lite and QFF-3D are better than MonoSDF with Multi-Resolution Grids, but worse than MonoSDF with MLP. Some artifacts appear in portions of the scene that are not well observed by the three views. We find our results consistent with those in MonoSDF: MLP performs better due to larger smoothness bias in sparse views, and the high-parameter models tend to overfit or produce spurious artifacts. Qualitatively, shown in Figure 6, our method can reconstruct thin regions, such as the arms and noses of the small snowman, which MLP cannot complete. MLP captures better geometry in smooth regions with high frequency textures, such as faces.

**Large scale scene reconstruction.** Similarly, we apply our QFF-Lite and QFF-3D on MonoSDF [40] for the ScanNet dataset [7]. We use  $N = 16$ ,  $M = 2^7$   $L = 6$ , 8-layer MLP for QFF-Lite and 2-layer MLP for QFF-3D for all scenes. Table 4 compares ScanNet evaluation to the baseline methods. Our QFF-3D achieves the best Chamfer- $L_1$  distance, and the best overall F-score, followed by our QFF-Lite. In more detail, we achieve significantly higher recall but slightly worse accuracy and precision compared to MLP and Multi-Res Grid. We also qualitatively verify the

Method	Acc. ↓	Comp. ↓	Chamfer- $L_1$ ↓	Prec. ↑	Recall ↑	F-score ↑	# Params
COLMAP [26]	0.047	0.235	0.141	0.711	0.441	0.537	-
UNISURF [18]	0.554	0.164	0.359	0.212	0.362	0.267	802 K
NeuS [33]	0.179	0.208	0.194	0.313	0.275	0.291	1.41 M
VolSDF [38]	0.414	0.120	0.267	0.321	0.394	0.346	802 K
Manhattan-SDF [10]	0.072	0.068	0.070	0.621	0.586	0.602	1.06 M
NeuRIS [32]	0.050	0.049	0.050	0.717	0.669	0.692	1.41 M
MonoSDF (Multi-Res. Grids) [40]	0.072	0.057	0.064	0.770	0.601	0.626	12.5 M
MonoSDF (MLP) [40]	<b>0.035</b>	0.048	0.042	<b>0.799</b>	0.681	0.733	<b>711 K</b>
MonoSDF (QFF-Lite)	0.043	0.044	0.044	0.761	0.718	0.738	1.63 M
MonoSDF (QFF-3D)	0.040	<b>0.041</b>	<b>0.041</b>	0.765	<b>0.744</b>	<b>0.754</b>	9.97 M

Table 4. **Evaluation on ScanNet [7].** We report the baseline and our results on ScanNet [7]. We mark bold for the best methods for each criteria. We use the results provided by [40] for the baseline.

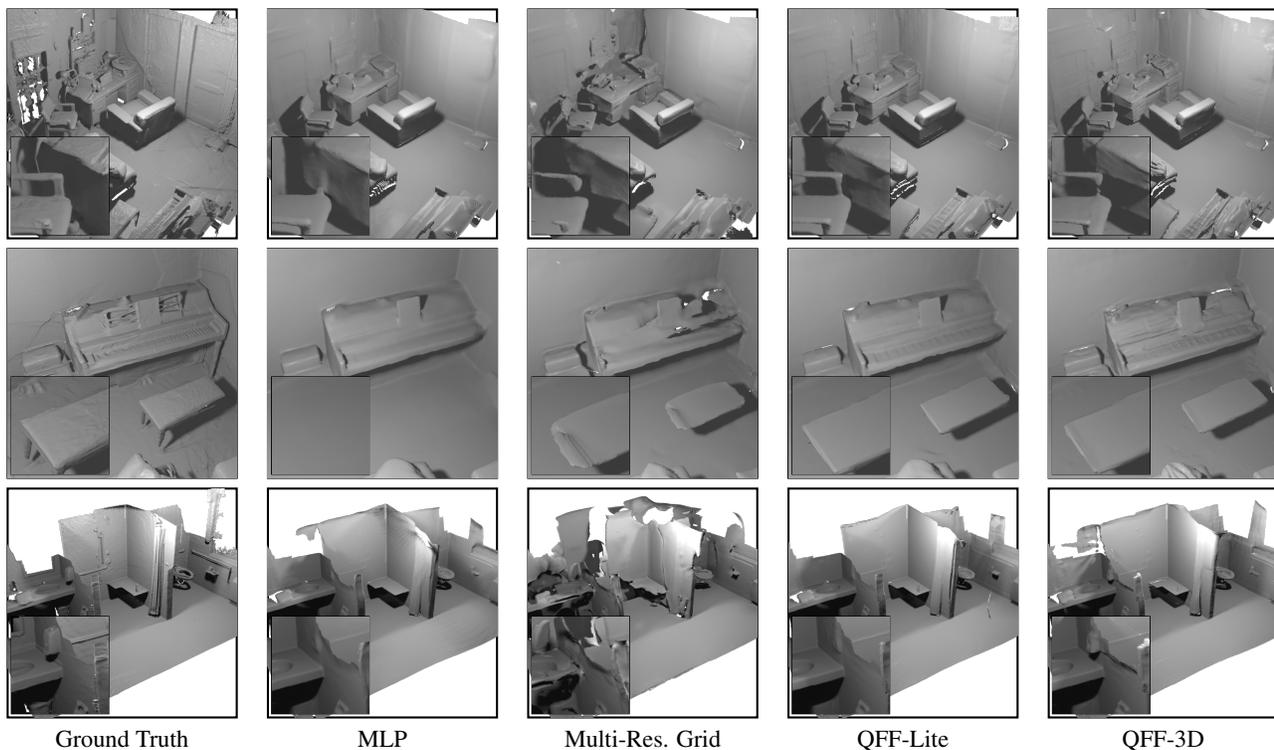


Figure 7. **Qualitative visualization of ScanNet Reconstruction.** MLP and Multi-Res Grid refers to the MonoSDF [40] with corresponding architectures. We use author provided meshes for visualizing the MLP, Multi-Res Grid. Best viewed in zoomed.

robustness of our method as shown in Figure 7. Our QFF-Lite and QFF-3D are capable of capturing correct geometry compared to Multi-Res Grid, and are on-par with MLP. We also emphasize that both of our methods are capable of capturing high-frequency geometries, such as piano keyboards or drawer hands.

#### 5.4. Ablation Studies and Discussions

We compare the impact of different design choices of our method on the *Lego* scenario of the Nerf Synthetic

Dataset [16]. We use the default values of both QFF-Lite and QFF-3D with  $N = 16$ ,  $M = 2^7$  and  $L = 6$ , and use 8-layer MLP for QFF-Lite and 2-layer MLP for QFF-3D. Table 5 shows a summary of our ablation studies.

**Resolution of Feature Bins.** Given the same feature vector length and number of layers, for QFF-Lite, we find a small increment of PSNR as the number of bins increases, but the increment slows down as we further increase the resolution, with 0.24 PSNR improvement from  $M = 2^5 \rightarrow 2^7$  and 0.04 PSNR from  $M = 2^7 \rightarrow 2^9$ . For QFF-3D, increasing

Method	Feats. (N)	Bins (M)	Layers	PSNR	# Params
QFF-Lite	8	$2^7$	8	35.41	786 K
	<b>16</b>	<b><math>2^7</math></b>	<b>8</b>	35.52	934 K
	32	$2^7$	8	35.83	1.45 M
	16	$2^5$	8	35.28	897 K
	16	$2^9$	8	35.56	1.45 M
	16	$2^7$	6	35.56	876 K
QFF-3D	8	$2^7$	2	36.22	5.00 M
	<b>16</b>	<b><math>2^7</math></b>	<b>2</b>	36.68	9.82 M
	32	$2^7$	2	36.91	19.5 M
	16	$2^5$	2	35.28	925 K
	16	$2^9$	2	36.00	151 M
	16	$2^7$	3	36.83	9.89 M
	16	$2^7$	4	36.82	9.96 M

Table 5. Comparison of different hyper-parameters. Bold values denote our default model. *Italicized* values denote the changes to the default model.

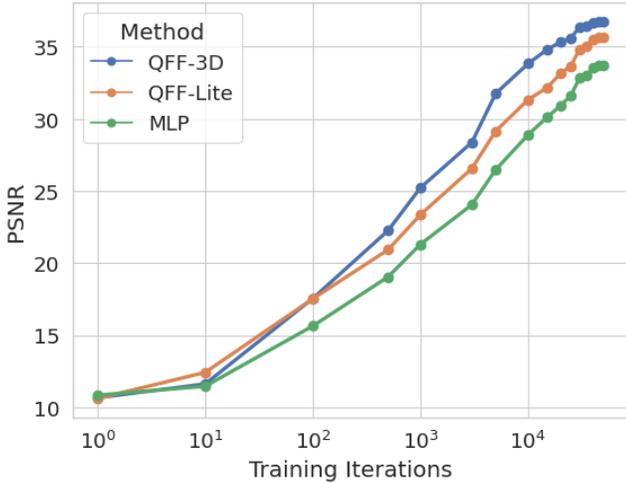


Figure 8. Comparison of Convergence for the Lego scene. We visualize number of steps taken vs. Test Image PSNR for our QFF-Lite and QFF-3D against MLP.

the number of bins to an excessively high number reduces the PSNR, because each bin in 2D is too fine-grained to be accessed multiple times. The number of parameters is proportional to  $\mathcal{O}(M)$  for QFF-Lite and  $\mathcal{O}(M^2)$  for QFF-3D.

**Length of Feature Vector.** We find that increasing the length of the feature vectors benefits both QFF-Lite and QFF-3D, improving 0.11, 0.46 PSNRs from  $N = 8 \rightarrow 16$ , and 0.31, 0.23 PSNRs from  $N = 16 \rightarrow 32$ , for QFF-Lite and QFF-3D respectively. Changing  $N$  changes the number of parameters proportional to  $\mathcal{O}(M)$  for QFF-Lite and  $\mathcal{O}(M^2)$  for QFF-3D.

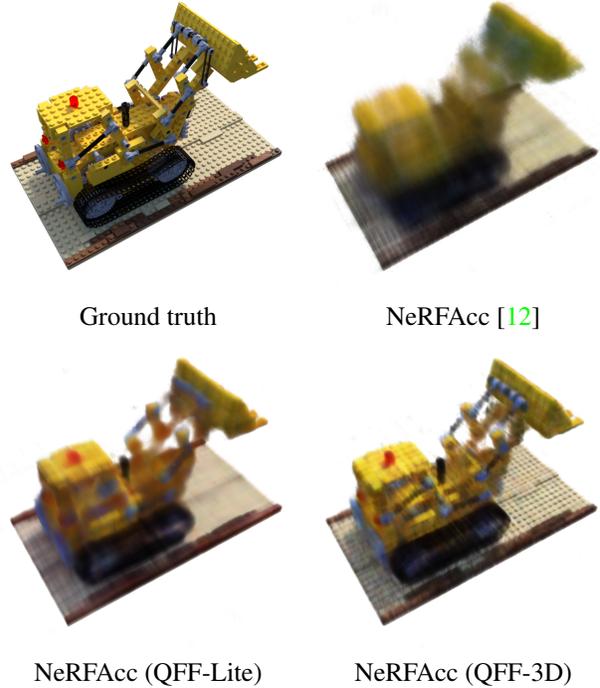


Figure 9. Visualization of the Lego scene at 500 steps of training. Best viewed in zoomed.

**Number of Layers.** We vary the number of layers in QFF-3D and find a slight improvement of 0.15 PSNR with MLP going from 2 to 3 layers, but 4-layer does not improve further.

**Convergence time compared to MLP.** In addition, we compare the convergence speed of the test-time loss between our method and the baseline 8-layer MLP. Figure 8 shows the plot of the convergence time of ours and the MLP. We show that our QFF-3D trains the fastest, followed by QFF-Lite, then MLP. We emphasize that the first 5000 iterations of our QFF-3D roughly correspond to 15000 iterations of QFF-Lite and 25000 iterations of MLP. Figure 9 visualizes the MLP, QFF-Lite and QFF-3D only after 500 iterations of training. We show that both our QFF-Lite and QFF-3D produce sharper renderings over MLP, and that QFF-3D generates high-frequency details within a few iterations.

## 6. Conclusion

We present Quantized Fourier Features (QFF), an easy-to-train, memory-efficient yet expressive neural field representation. QFF combines the best worlds of explicit feature representation and frequency-based positional encoding. We demonstrate advantages of QFF in wide range of applications of neural field representations.

**Acknowledgements.** This research is partially supported by NSF IIS 2020227, a gift from Amazon and a gift from Illinois-Inspire Collaborative Research Fund.

## References

- [1] Henrik Aanaes, Rasmus Ramsbøl Jensen, George Vogiatzis, Engin Tola, and Anders Bjorholm Dahl. Large-scale data for multiple-view stereopsis. *International Journal of Computer Vision (IJCV)*, 2016. 1, 6
- [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 *IEEE International Conference on Computer Vision (ICCV)*, 2021. 2
- [3] Steven L Brunton, Bernd R Noack, and Petros Koumoutsakos. Machine learning for fluid mechanics. *Annual review of fluid mechanics*, 2020. 2
- [4] Rohan Chabra, Jan Eric Lenssen, Eddy Ilg, Tanner Schmidt, Julian Straub, Steven Lovegrove, and Richard Newcombe. Deep local shapes: Learning local sdf priors for detailed 3d reconstruction. In *European Conference on Computer Vision (ECCV)*, 2020. 2
- [5] Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J Guibas, Jonathan Tremblay, Sameh Khamis, et al. Efficient geometry-aware 3d generative adversarial networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [6] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. In *European Conference on Computer Vision (ECCV)*, 2022. 2, 3, 4, 5, 6
- [7] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 6, 7
- [8] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 3
- [9] Amos Gropp, Lior Yariv, Niv Haim, Matan Atzmon, and Yaron Lipman. Implicit geometric regularization for learning shapes. In *ICML*, 2020. 2
- [10] 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 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 6, 7
- [11] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 4, 5
- [12] Ruilong Li, Matthew Tancik, and Angjoo Kanazawa. Nerfacc: A general nerf acceleration toolbox. *arXiv preprint arXiv:2210.04847*, 2022. 5, 8
- [13] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 2
- [14] Julien Martel, David Lindell, Connor Lin, Eric Chan, Marco Monteiro, and Gordon Wetzstein. Acorn: Adaptive coordinate networks for neural scene representation. *ACM Transactions on Graphics (TOG)*, 2021. 2
- [15] Ricardo Martin-Brualla, Noha Radwan, Mehdi S. M. Sajjadi, Jonathan T. Barron, Alexey Dosovitskiy, and Daniel Duckworth. NeRF in the Wild: Neural Radiance Fields for Unconstrained Photo Collections. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [16] 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. 1, 2, 3, 4, 5, 7
- [17] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (TOG)*, 2022. 1, 2, 3, 5
- [18] Michael Oechsle, Songyou Peng, and Andreas Geiger. Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction. In *IEEE International Conference on Computer Vision (ICCV)*, 2021. 7
- [19] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 1, 2
- [20] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems (NeurIPS)*, 2019. 3
- [21] Songyou Peng, Michael Niemeyer, Lars Mescheder, Marc Pollefeys, and Andreas Geiger. Convolutional occupancy networks. In *European Conference on Computer Vision (ECCV)*, 2020. 3
- [22] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 3
- [23] Nasim Rahaman, Aristide Baratin, Devansh Arpit, Felix Draxler, Min Lin, Fred Hamprecht, Yoshua Bengio, and Aaron Courville. On the spectral bias of neural networks. In *International Conference on Machine Learning (ICML)*. PMLR, 2019. 1, 2
- [24] Maziar Raissi, Paris Perdikaris, and George E Karniadakis. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational physics*, 378:686–707, 2019. 1, 2
- [25] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. In *IEEE International Conference on Computer Vision (ICCV)*, 2021. 1, 2
- [26] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise View Selection for Unstructured Multi-View Stereo. In *European Conference on Computer Vision (ECCV)*, 2016. 6, 7
- [27] Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 2, 3, 4

- [28] Pratul P Srinivasan, Boyang Deng, Xiuming Zhang, Matthew Tancik, Ben Mildenhall, and Jonathan T Barron. Nerv: Neural reflectance and visibility fields for relighting and view synthesis. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [29] Yannick Strümpfer, Janis Postels, Ren Yang, Luc Van Gool, and Federico Tombari. Implicit neural representations for image compression. In *European Conference on Computer Vision (ECCV)*, 2022. 1
- [30] Towaki Takikawa, Joey Litalien, Kangxue Yin, Karsten Kreis, Charles Loop, Derek Nowrouzezahrai, Alec Jacobson, Morgan McGuire, and Sanja Fidler. Neural geometric level of detail: Real-time rendering with implicit 3d shapes. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [31] Matthew Tancik, Pratul P. Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan T. Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 1, 2, 4
- [32] Jiepeng Wang, Peng Wang, Xiaoxiao Long, Christian Theobalt, Taku Komura, Lingjie Liu, and Wenping Wang. Neuris: Neural reconstruction of indoor scenes using normal priors. In *European Conference on Computer Vision (ECCV)*, 2022. 7
- [33] 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 (NeurIPS)*, 2021. 3, 7
- [34] Suttisak Wizadwongsa, Pakkapon Phongthawee, Jiraphon Yenphraphai, and Supasorn Suwajanakorn. Nex: Real-time view synthesis with neural basis expansion. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [35] Liwen Wu, Jae Yong Lee, Anand Bhattad, Yuxiong Wang, and David Forsyth. Diver: Real-time and accurate neural radiance fields with deterministic integration for volume rendering. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16179–16188, 2022. 2, 3
- [36] Zhiqin Xu. Frequency principle: Fourier analysis sheds light on deep neural networks. *Communications in Computational Physics*, 2020. 1, 2
- [37] Zhi-Qin John Xu, Yaoyu Zhang, and Tao Luo. Overview frequency principle/spectral bias in deep learning. *arXiv preprint arXiv:2201.07395*, 2022. 2
- [38] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces. *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. 7
- [39] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenotrees for real-time rendering of neural radiance fields. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [40] Zehao Yu, Songyou Peng, Michael Niemeyer, Torsten Sattler, and Andreas Geiger. Monosdf: Exploring monocular geometric cues for neural implicit surface reconstruction. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. 1, 2, 6, 7
- [41] Xiaoshuai Zhang, Sai Bi, Kalyan Sunkavalli, Hao Su, and Zexiang Xu. Nerfusion: Fusing radiance fields for large-scale scene reconstruction. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 3