

TriNeRFLet: A Wavelet Based Multiscale Triplane NeRF Representation

Rajaei Khatib
 Electrical Engineering
 Tel Aviv University
 rajaekh@mail.tau.ac.il

Raja Giryes
 Electrical Engineering
 Tel Aviv University
 raja@tauex.tau.ac.il

<https://rajaekh.github.io/trinerflet-web>

Abstract

In recent years, the neural radiance field (NeRF) model has gained popularity due to its ability to recover complex 3D scenes. Following its success, many approaches proposed different NeRF representations in order to further improve both runtime and performance. One such example is Triplane, in which NeRF is represented using three 2D feature planes. This enables easily using existing 2D neural networks in this framework, e.g., to generate the three planes. Despite its advantage, the triplane representation lagged behind in its 3D recovery quality compared to NeRF solutions. In this work, we propose TriNeRFLet, a 2D wavelet-based multiscale triplane representation for NeRF, which closes the 3D recovery performance gap and is competitive with current state-of-the-art methods. Building upon the triplane framework, we also propose a novel super-resolution (SR) technique that combines a diffusion model with TriNeRFLet for improving NeRF resolution.

1. Introduction

3D scene reconstruction from multiple 2D views is a challenging task that has been widely studied, and many methods have been proposed to solve it. Neural radiance field (NeRF) [27] is prominent among these methods as it has demonstrated a high level of generalizability in generating novel views with high quality and consistent lighting. NeRF utilizes an implicit representation of the 3D scene in the form of a multi-layer perceptron (MLP). This enables it to capture complex 3D geometry and lighting.

At a high level, NeRF relies on the rendering equation that approximates each pixel in the image using points sampled along the ray that passes through it [27]. The MLP in NeRF takes as input the frequency encoding of the Euclidean coordinates and view direction of the point of interest, and outputs the radiance and density at this point.

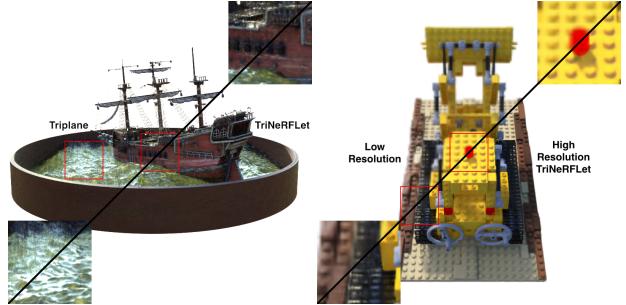


Figure 1. Our approach improves the quality of NeRF reconstruction. (Left) Using a multiscale wavelet representation with Triplane improves 3D reconstruction quality; (Right) Combining diffusion-based super-resolution with TriNeRFLet leads to high-resolution reconstruction from low-resolution 2D views.

The MLP weights are learned via end-to-end optimization, by comparing the values between the rendered pixel and its value in the corresponding 2D image. After the training process is over, the MLP can be used to render novel views.

Due to its success in representing 3D scenes, many methods proposed improvements to NeRF [5, 6, 28]. They sought to enhance NeRF 3D reconstruction capabilities, as well as other drawbacks such as high runtime and aliasing artifacts that it suffered from. One such approach uses three axis-aligned 2D feature planes, denoted as Triplane, to represent the NeRF [7]. In the rendering process, each point is sampled by projecting it onto each of the three planes and then concatenating the features that correspond to the three projections. This forms a single feature vector for the point that is then passed to a small MLP that outputs the density and color values of this point.

A significant advantage of the Triplane representation is that it can be used with many already existing 2D methods. In the original work [7], the authors used an existing 2D Generating Adversarial Network (GAN) architecture to generate its planes. Follow-up works employed the 2D

property of the Triplane to perform NeRF super-resolution [4] and 3D generation [21, 37].

While being useful due to its special 2D structure, the reconstruction quality achieved by Triplane lagged behind other efficient multiview reconstruction methods such as instang NGP (INGP) [28], which is an improvement of NeRF, and 3D Gaussian splatting [20]. Due to its structure, only Triplane entries that are included in the training views rays are learned. Thus, there might be entries at the planes that will still have their initial random values also after training finishes. These random features will be used by novel views that use them and therefore will deteriorate the quality of their generation and create artifacts in them.

To tackle these drawbacks, we present a novel Triplane representation that relies on a multiscale 2D wavelet structure. Instead of learning the feature planes directly, we learn the wavelet features of several resolutions, while regularizing the wavelet representation to be sparse. Thus, regions that are covered by the training views are learned in a similar way to the case of a regular Triplane, and regions that are not covered by the training views are updated by a lower resolution estimate according to how much they are affected by their neighboring regions. Figure 1(left) shows the improvement in reconstruction quality attained by this representation. Figure 2 illustrates how regions on the plane may affect each other due to the wavelet representation.

We use L_1 regularization on the wavelet coefficients. The goal is to sparsify the coefficients that are not trained by the given views such that only the coefficients that were trained will be used. Our proposed structure is inspired by the wavelet literature, which indicates that wavelet features can represent data accurately while being sparse [11].

To further improve the training, we also perform the training in a multi-scale fashion, starting with a lower resolution version of the images that updates only the coarse coefficients of the wavelet. Then, after some iterations we increase back the resolution and add more resolution to the wavelet representation.

Our proposed TriNeRFLet solution closes the performance gap of Triplane against other multiview 3D reconstruction methods and makes it competitive with current state-of-the-art (SOTA) methods. In addition, we show how we may use pre-trained diffusion models for 2D super-resolution (SR) [33, 34] with our TriNeRFLet representation in order to propose a novel NeRF SR approach. Figure 1(Right) demonstrates the improvement this approach attains. We show on various experiments the advantage of our TriNeRFLet in novel view reconstruction compared to INGP and 3D Gaussian splatting and the superiority of our SR solution compared to competing approaches that do not require dedicated training for multiview or 3D data.

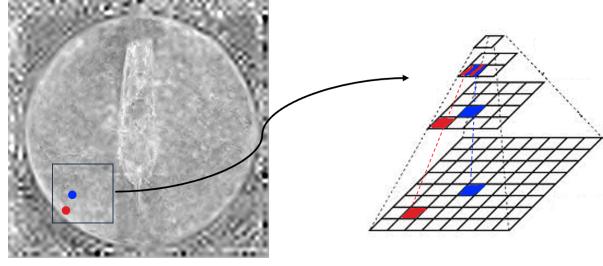


Figure 2. The multiscale wavelet representation transfers information between points on the triplane planes: Although the Red dot on the plane was not covered by the training images views, it will get a coarser estimate from a lower resolution wavelet layer that corresponds also to the blue dot, which was updated in training.

2. Related Works and Background

NeRF [27] is a framework for reconstructing a 3D scene from a set of multiview images. Its main component is an implicit neural representation of the 3D scene via a mapping neural network $F_\theta : (\mathbf{x}, \mathbf{d}) \rightarrow (c, \sigma)$, where \mathbf{x} is the Cartesian coordinate of the point of interest, \mathbf{d} is the direction vector from it to the camera origin, c is the RGB color at this point and σ is the density.

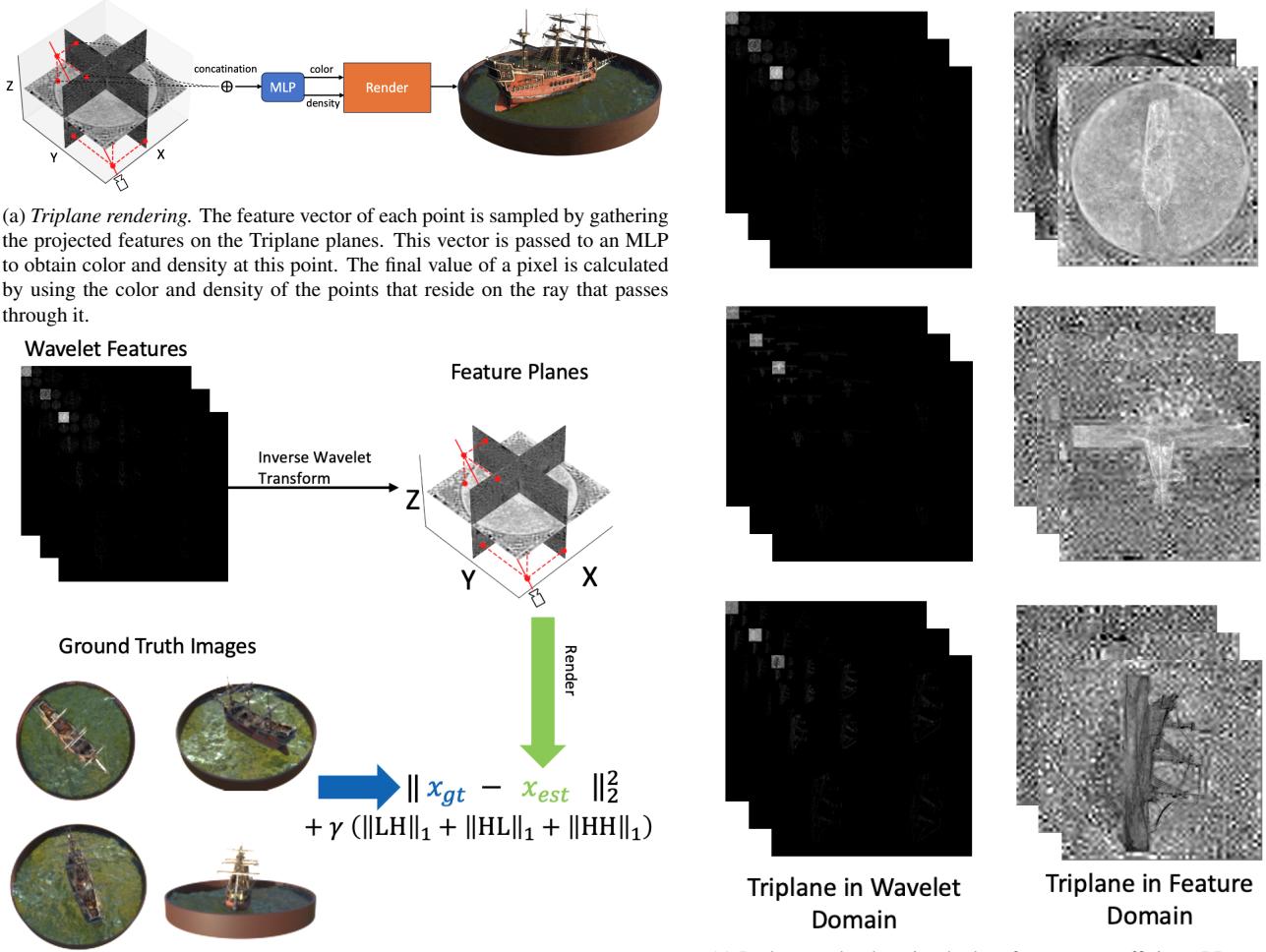
The network F is trained using the pixels of the multiview images using the rendering equation. For a given pixel, it relates its value to N points sampled along the ray that passes through this pixel, where the ray direction is derived from the image view. Denoting by $t_1 < t_2 < \dots < t_N$ the depth of the points, $\delta_i = t_{i+1} - t_i$ the distance between adjacent samples, and σ_i, c_i the density and color of each point respectively, then the pixel value \hat{C} can be estimated from these points via:

$$\hat{C} = \sum_{i=1}^N T_i (1 - \exp(-\sigma_i \delta_i)) c_i, \quad (1)$$

where $T_i = \exp \left(- \sum_{j=1}^{i-1} \sigma_j \delta_j \right).$

This equation is a discrete approximation of the rendering equation. As we have the true pixel value C of the image, the network parameters θ are trained simply by minimizing the L2 loss $\|C - \hat{C}\|_2$.

In the classical NeRF scheme [27], the mapping function F_θ is divided into two main parts. The first part is a frequency encoding function that encodes \mathbf{x} and \mathbf{d} using Fourier features [38]. These vectors are then concatenated to form a feature vector that is fed to the second part, which is an MLP function that maps the feature vector into its corresponding color and density values. Some extensions to NeRF consider optimization using unconstrained photo collections [25], few view points [42], non-rigidly deform-



(b) *TriNeRFLet training scheme*. First, wavelet features are transformed into the Triplane domain. Next, pixels are rendered using these features in order to fit them to their ground-truth values. The high frequencies channels of LH, HL and HH from all wavelet levels get regularized by L_1 loss.

Figure 3. The *TriNeRFLet* framework learns features in wavelet domain, which are transformed into feature domain to render 3D objects.

ing scenes [29], and imperfect camera poses [22].

INGP [28] is an important development, which managed to decrease NeRF runtime significantly while getting even better reconstruction performance. It achieved that by using a multiscale hash grid structure alongside a small shallow MLP that represents the scene instead of the big deep one used in previous works. Furthermore, this method introduced new rendering techniques and low-level optimizations in order to get a faster runtime. The multiscale concept of INGP will also be used in our proposed TriNeRFLet.

3D Gaussian Splatting [20] is a state-of-the-art 3D reconstruction approach that takes a different strategy than NeRF. Instead of using an implicit neural function, it represents the scene as a group of 3D Gaussians, and learns their parameters. They introduce a new method for efficiently rendering these Gaussians, thus, managing to train

rapidly and achieve 3D rendering in real-time. In run-time, one needs to distinguish between train-time and rendering time. The train time of this method is similar to INGP but its rendering is faster than current NeRF solutions.

Triplane [4, 7] is a NeRF variant, in which the frequency encoding part is replaced by three perpendicular 2D feature planes of dimension $N \times N$ with C channels. To get the feature vector of a point x , the point is projected onto each of these planes, then the sampled feature values of the projected points are gathered together to form a feature vector. And the final feature vector is obtained by concatenating this feature vector with the frequency encoding of the direction vector d . These feature planes are learned alongside the MLP weights with a similar approach to what was explained earlier. This process is illustrated in Figure 3a.

The advantage of Triplane compared to other alternatives

is its use of 2D planes, which eases the use of standard 2D networks with it. This has been used in various applications. Bahat et al. [4] utilized it to train a NeRF super-resolution network from multiview training data. By exploiting the triplane structure, a 2D super resolution neural network was trained to perform super resolution on low resolution planes and in this way increase the overall NeRF resolution. Another example is [37], in which a diffusion model was trained to generate a Triplane, thus utilizing the high generalization power of diffusion models and bringing it to 3D generation. Similarly, Triplane is predicted by a transformer-based neural network in [21] in order to have fast text to 3D generation.

Wavelets were widely used in signal and image processing applications [24]. With the advent of deep learning, they were employed in various architectures to improve performance and efficiency. The filter-bank approach of wavelet was utilized to represent each kernel in a neural network as a combination of filters from a given basis, e.g., a wavelet basis [31]. The concept of using multi-resolution was used to improve positional encoding of implicit representations [15, 23, 36]. Wavelets were used to improve efficiency and generation of StyleGAN [10], diffusion models [13, 17, 30] and normalizing flows [43].

Rho *et al.* [32] implemented a multiscale 3D grid using a grid decomposition method that transforms 2D planes alongside 1D vectors (unrelated to the 2D planes) into a 3D grid, where they used a 2D multiscale wavelet representation to represent the 2D planes. Despite the use of 2D multiscale representation, this method is constrained by 3D grid limitations, namely, it consumes a lot of memory and is limited to low resolutions only. In this work, we deploy the 2D wavelet multiscale representation on Triplane planes directly, thus being freed from 3D grid limitations with a memory-efficient and high-resolution representation.

Diffusion Models [9, 16] have become a very popular tool in recent years for image manipulation and reconstruction tasks [18, 19, 33, 35, 41], where a denoising network is trained to learn the prior distribution of the data. Diffusing models can be used also beyond generation, where given a degraded image, some conditioning mechanism is combined with the learned prior to solve different tasks [2, 3, 8]. For example, in [1, 35, 41] diffusion models were utilized for the problems of deblurring and super-resolution. In this work, we use the Stable Diffusion (SD) upscaler [33] to perform super-resolution on 2D projections of the low-resolution NeRF in order to create a high-resolution version of the NeRF. While the SD upscaler is conditioned both on the low-resolution image and a given text prompt, we use it without text, as explained later in the paper.

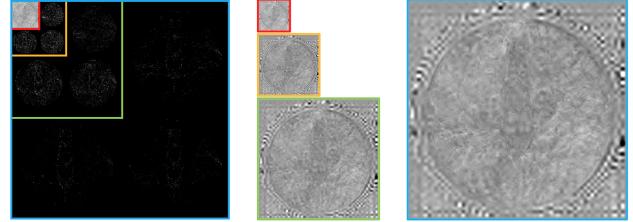


Figure 4. Example of the multiscale property of the wavelet representation. On the right, we see the same feature plane at different resolutions. On the left we see the wavelet representation of the plane at the left, which is of the highest resolution. Note that its wavelet representation includes the wavelet representations of the lower-resolution planes (the colors correspond to the matching between each plane and its wavelet representation).

3. TriNeRFLet

To alleviate the Triplane drawbacks mentioned above, we present TriNeRFLet. Instead of using the 2D feature planes directly, we optimize them in the wavelet domain. At a high level, wavelet transforms a 2D plane to a multiscale representation that contains both low and high frequencies at different resolutions. The inverse wavelet transform converts this representation back to the original 2D plane. Note that we apply the wavelet separately on each channel. We use the following notation for the wavelet representation. At resolution i of the wavelet representation, we denote by LL_i the low frequencies of this resolution, and by LH_i , HL_i and HH_i the high frequencies. Note that LL_i is used to create the coefficients of the next resolution of the wavelet representation. At a certain step, we stop and remain with LL .

In rendering, TriNeRFLet is similar to regular Triplane, in the sense that given a point it calculates its features by projecting it to the Triplane and then feed the features to a MLP that outputs the color c_i and density σ_i of the point that are then used by the rendering equation (see Figure 3a). Yet, instead of holding the Triplane in the feature domain, TriNeRFLet holds it in the Wavelet domain. Thus, in the optimization, instead of optimizing the values of the triplane directly, we optimize the wavelet coefficients of each pixel (see Figure 3b). Note that in the regular Triplane scheme, each plane entry impacts only one location in each plane. In TriNeRFLet, it impacts different resolutions in the wavelet representation of the plane (see Figure 2). Figure 3c shows some learned planes for the ship object and their corresponding Wavelet coefficients.

The learned parameters of TriNeRFLet, alongside the MLP weights, are the wavelet coefficients, and they are learned in an end-to-end fashion. The multiscale structure of the wavelet allows the method to learn fine details whenever it needs, and to have a lower resolution estimate from lower resolution wavelet layers that are affected by nearby values when the area is not covered by the training pixels.

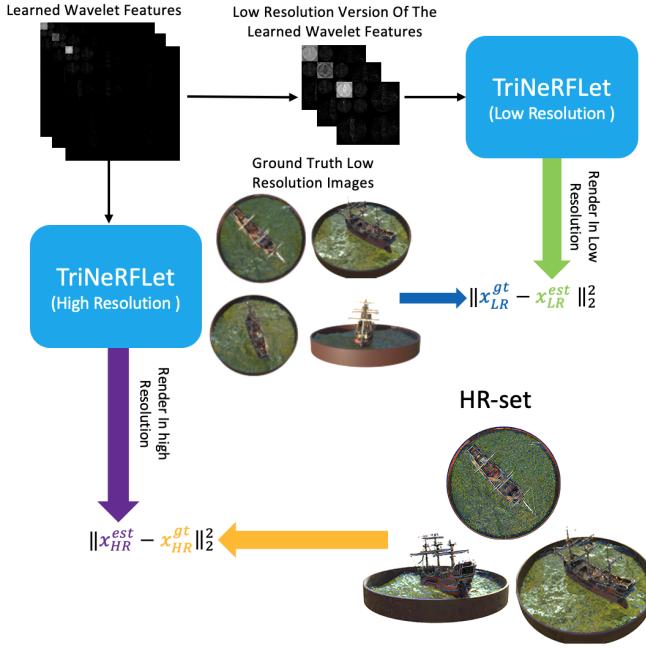


Figure 5. *TriNeRFLet super-resolution*. Low-resolution *TriNeRFLet* renders LR images using the LR version of the wavelet features, and fits them to the given low-resolution images. High-resolution *TriNeRFLet* renders HR images using all wavelet levels, and fits them to HR images from HR-set.

The question remains, what should be done with the coefficients of the higher frequencies that are not updated. From wavelet theory [11], we expect the LH, HL and HH channels to be sparse. Thus, we regularize these channels to be sparse over all resolutions by adding an L_1 regularization for these coefficients to the training loss. This allows us to increase the resolution of the Triplane without having the problem of having features in it that remain random. In *TriNeRFLet*, also these features get updated. The overall loss of *TriNeRFLet* is:

$$\begin{aligned} \text{loss} = & \sum_{i \in \text{TrainingPixels}} \left\| C_i - \hat{C}_i \right\|_2^2 + \\ & \gamma \sum_{l \in \text{resolutions}} (\|LH_l\|_1 + \|HL_l\|_1 + \|HH_l\|_1), \end{aligned} \quad (2)$$

where γ is the L_1 regularization factor.

To further comprehend the significance of this multiscale wavelet structure against regular Triplane, we compare the behaviour of each method when rendering novel views. In regular Triplane, some features will still have their initial random values after training ends. This introduces artifacts when rendering novel views, as demonstrated in the ship example in Figure 1. Note that these artifacts are absent in *TriNeRFLet* because all the features in the Triplane are updated either by direct training or using a coarser estimate from lower resolution wavelet coefficients (see Figure 2).

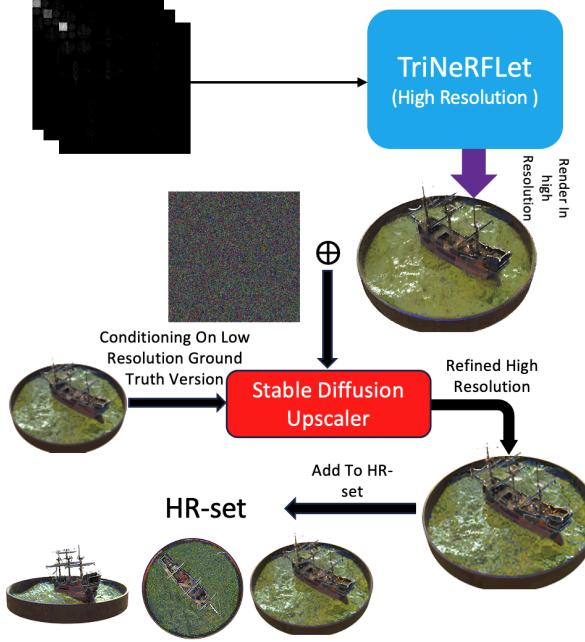


Figure 6. *SD_{refine}* process. First, a HR image is rendered. To improve its quality, noise is added to it and then it is plugged into Stable Diffusion upscaler with its LR version for conditioning. The result, which is a refined HR image is added to the HR-set.

This multiscale wavelet approach enjoys many advantages. First, the L_1 regularization on the wavelet coefficients prevents the method from over-fitting. In regular NeRF, INGP or Triplane, an over-fitting occurs when training for too many iterations. We have found empirically that this does not happen in our approach.

Another advantage is that the multiscale property of the wavelet representation enables increasing the resolution of the *TriNeRFLet* during training seamlessly. It is possible to continue to train the higher-resolution *TriNeRFLet* while using the wavelet coefficients of the lower-resolution *TriNeRFLet*. The wavelet of the lower-resolution plane is a subset of the wavelet of the higher resolution plane (see Figure 4). We use this flexibility to improve the training by learning the wavelet planes in a coarse to fine manner.

The multiscale wavelet structure is not ideal in terms of training runtime due to the fact that the feature planes need to be reconstructed in each training step. Yet, its rendering runtime is quite fast since the feature planes are reconstructed once only and then the runtime is competitive with other well-known fast NeRF schemes like INGP [28].

4. *TriNeRFLet Super-Resolution*

Given our *TriNeRFLet* scheme, we present a novel NeRF super-resolution approach that relies on the *TriNeRFLet* special multiscale structure. This method employs a 2D pre-trained stable diffusion (SD) upscaler [33] for super-

Algorithm 1 TriNeRFLet Super-Resolution

Input: Ground truth low resolution (LR) images $\{\mathbf{x}_{i,LR}^{gt}\}_{i=1}^k$, LR only steps s_{low_res} , Total steps s , Wavelet LR levels L_{LR} , Refresh steps $s_{refresh}$, Diffusion step limits T_{min}, T_{max} .

Output: High-Resolution (HR) TriNeRFLet

```

1: HR-set  $\leftarrow \{\}$                                  $\triangleright$  empty dictionary
2: for  $i = 0 \dots s - 1$  do
3:   Choose randomly a frame i -  $\mathbf{x}_{i,LR}^{gt}$ .
4:   Render  $\mathbf{x}_{i,LR}^{est}$  using only first  $L_{LR}$  wavelet levels.
5:    $loss \leftarrow \|\mathbf{x}_{i,LR}^{est} - \mathbf{x}_{i,LR}^{gt}\|_2^2$ 
6:   if  $i \geq s_{LR}$  then
7:     Render  $\mathbf{x}_{i,HR}^{est}$  using all wavelet levels.
8:
9:     if  $i \% s_{refresh} == 0$  then
10:      HR-set  $\leftarrow \{\}$                                  $\triangleright$  empty dictionary
11:      end if
12:
13:      if  $i \in \text{HR-set}$  then
14:         $\mathbf{x}_{i,HR}^{gt} \leftarrow \text{HR-set}[i]$ 
15:      else
16:         $t \sim \text{Rand}(T_{min}, T_{max})$ 
17:         $\mathbf{x}_{i,HR}^{gt} \leftarrow \text{SD}_{\text{refine}}(\mathbf{x}_{i,HR}^{est}, \mathbf{x}_{i,LR}^{gt}, t)$ 
18:        HR-set[i]  $\leftarrow \mathbf{x}_{i,HR}^{gt}$ 
19:      end if
20:       $loss \leftarrow loss + \|\mathbf{x}_{i,HR}^{est} - \mathbf{x}_{i,HR}^{gt}\|_1 + \lambda \text{LPIPS}(\mathbf{x}_{i,HR}^{est}, \mathbf{x}_{i,LR}^{gt})$ 
21:    end if
22:     $loss \leftarrow loss + \gamma \sum_l (\|\text{LH}_l\|_1 + \|\text{HL}_l\|_1 + \|\text{HH}_l\|_1)$ 
23:     $T_{max} \leftarrow \text{scheduler}(T_{max}, i)$                  $\triangleright$  decrease  $T_{max}$ 
24:  end for

```

resolution "guidance". Our approach does not require any kind of low-high resolution 3D scene pairs for supervision as done in other works [4]. Utilizing the robustness of the SD upscaler, our approach has the potential to handle different types of scenes and a range of resolutions.

Given a set of low-resolution (LR) multiview images of a given scene, the goal is to generate high-resolution novel views for this scene. We fit a TriNeRFLet to the LR images and aim at estimating a high-resolution (HR) TriNeRFLet. Denote the size of the LR planes by N_{LR} and the number of wavelet levels by L_{LR} . For the HR planes, denote the size by N and the number of levels by L . Due to the multiscale property of the wavelet representation, in the HR TriNeRFLet we just need to learn the high frequencies of the wavelet representation, i.e., the first L_{LR} levels of the HR wavelet planes are shared with the LR TriNeRFLet.

Figure 5 illustrates our TriNeRFLet super-resolution strategy. First, we start by fitting a LR TriNeRFLet to the provided low-resolution multiview images. After training

	N_{LL}	L	N_{base}	N_{final}	C	γ	W	D_{dens}	D_{col}	train steps
Small	64	4	512	1024	16	0.2	64	1	2	6k
Base	64	5	512	2048	32	0.4	64	1	2	10k
Light										
Base	64	5	512	2048	32	0.4	64	1	2	43k
Large	64	5	512	2048	48	0.6	128	1	2	83

Table 1. The parameters used in the different TriNeRFLet versions

for s_{LR} steps, the super-resolution process starts, and a key component of it is the stable diffusion upscaler refinement ($\text{SD}_{\text{refine}}$) step presented in Figure 6. Given a low resolution ground truth image \mathbf{x}_{LR}^{gt} , its high resolution version \mathbf{x}_{HR}^{est} is rendered using all wavelet levels L using the TriNeRFLet. Initially, the result will not be of high-resolution. Thus, to improve it we use a diffusion step with time t randomly selected from the range $[T_{min}, T_{max}]$. Then, similar to [14], a noise with variance that depends on t is added to \mathbf{x}_{HR}^{est} . Then we plug the noisy version of \mathbf{x}_{HR}^{est} into the SD upscaler to perform the diffusion process from step t until the end, while being conditioned on the LR image \mathbf{x}_{LR}^{gt} . This results is the enhanced version of \mathbf{x}_{HR}^{est} , denoted as $\mathbf{x}_{HR}^{enhanced}$. We denote this enhancement process by $\text{SD}_{\text{refine}}$.

We use the SD upscaler with a small modification compared to the original work [33]. In their work, the denoising step is $\mathbf{x}_{t-1} = \epsilon(\mathbf{x}_t, \mathbf{z}, \emptyset) + \alpha(\epsilon(\mathbf{x}_t, \mathbf{z}, \mathbf{y}) - \epsilon(\mathbf{x}_t, \mathbf{z}, \emptyset))$, where ϵ , α , \mathbf{z} and \mathbf{y} are the U-net denoiser, a guidance factor, the low-resolution image and the guidance text respectively. Instead of doing the guidance in the text direction, we do it in the low-resolution image direction and without text, i.e., $\mathbf{x}_{t-1} = \epsilon(\mathbf{x}_t, \emptyset, \emptyset) + \alpha(\epsilon(\mathbf{x}_t, \mathbf{z}, \emptyset) - \epsilon(\mathbf{x}_t, \mathbf{z}, \emptyset))$.

Having the refined HR images, we want to use them to update the TriNeRFLet. We add these images to the set denoted HR-set. This set is initialized to be empty and gets updated sequentially. It is refreshed (i.e. emptied) every $s_{refresh}$ steps as illustrated in Algorithm 1. The images in HR-set are used as "ground-truth" images for training the HR TriNeRFLet to update all its L levels.

To tie all threads together, the method starts by fitting TriNeRFLet with low-resolution features to the low-resolution ground-truth images. Then $\text{SD}_{\text{refine}}$ is used to improve the resolution of the TriNeRFLet generated views and the images are stored in HR-set. Next, a TriNeRFLet with high-resolution features (using all wavelet levels) is fitted to HR-set and simultaneously the low-resolution images are used to train the low levels of the wavelet representation of the same TriNeRFLet. As in regular TriNeRFLet training, a L_1 regularization of the wavelet coefficients is added to the loss. In addition, we use also a LPIPS loss between the rendered high resolution and the ground truth low resolution. As mentioned earlier, HR-set gets refreshed every $s_{refresh}$

Method	Mic	Chair	Ship	Materials	Lego	Drums	Ficus	Hotdog	Avg.	Train Time ↓	Render FPS ↑
Vanilla NeRF	32.91	33.00	28.65	29.62	32.54	25.01	30.13	36.18	31.005	≥ 1 day	≤ 0.2
INGP	36.22	35.00	31.10	29.18	36.39	26.02	33.51	37.40	33.176	~ 5 mins	0.6 – 1.2
3D Gaussian Splatting	35.36	<u>35.83</u>	30.80	<u>30.00</u>	35.78	<u>26.15</u>	34.87	<u>37.72</u>	33.32	~ 10 mins	135
Triplane	33.85	32.83	29.58	28.15	34.70	24.86	30.35	35.80	31.26	60-180 mins	1 – 4
Ours: TriNeRFLet Small	35.18	33.76	30.33	29.14	35.32	25.66	33.34	36.24	32.37	12-15 mins	1 – 4
Ours: TriNeRFLet Base Light	35.77	35.00	31.10	29.35	36.44	25.98	33.96	36.93	33.07	60-90 mins	1 – 4
Ours: TriNeRFLet Base	<u>36.72</u>	35.54	<u>31.77</u>	29.72	<u>36.66</u>	26	33.7	37.31	<u>33.43</u>	5-8 hours	1 – 4
Ours: TriNeRFLet Large	37.22	36	32.7	30.34	37.32	26.2	<u>34.47</u>	37.89	34.017	~ 1 day	1 – 4

Table 2. Blender Dataset PSNR (↑) results with train time and render FPS for each method. **Bold** is best, underline is second.

training steps. Furthermore, the diffusion step upper limit T_{max} is scheduled to linearly decrease during training in order to force SD_{refine} to introduce fewer changes as training converges. This is described in Algorithm 1.

The TriNeRFLet multiscale wavelet structure plays a crucial role in our SR scheme. First, its multiscale structure establishes a profound connection between high-resolution and low-resolution renderings by forcing them to share information (using the same first L_{LR} wavelet levels). Second, wavelet regularization, especially in the layers that only the HR TriNeRFLet uses, constrains the optimization process and forces the HR TriNeRFLet to learn only useful wavelet coefficients that contribute to HR details. Thus, this unique combination of multiscale wavelet and diffusion models allows the method to learn the high-resolution details it exactly needs without redundancy and artifacts.

The diffusion SR model that we use only supports SR from 128 to 512. However, we want to be able to work with other resolutions, such as 100 to 400 or 200 to 800. For lower resolutions, e.g., going from 100 to 400, we pad x_{LR}^{gt} and x_{HR}^{est} with zeros in order to achieve the desired resolution before plugging them into the SD upscaler. When it finishes, the result is un-padded to the original resolution. For higher resolutions, e.g., going from 200 to 800, we randomly crop 128 and 512 resolution crops from x_{LR}^{gt} and x_{HR}^{est} , respectively, that correspond to the same relative location. We plug the images into the upscaler, and the result is actually an enhanced version of the crop from x_{HR}^{est} . Then we only fit the high-resolution render to this crop instead of the whole image. When the HR-set is refreshed, a different random crop will be chosen, thus covering more areas as the process continues.

5. Experiments

We turn to present the results of our method. More visual results appear in the [project page](#) (preferred to open in Chrome browser). An ablation study appears in Appendix A.

5.1. 3D Scene Reconstruction

For evaluating TriNeRFLet in the 3D recovery task, we define four versions of TriNeRFLet - small, base light, base and large. The parameters of TriNeRFLet are wavelet LL component resolution - N_{LL} , base resolution - N_{base} , wavelet layers - L , final resolution - N_{final} , number of channels - Ch , wavelet regularization γ , MLP dim (neurons) - W , MLP hidden layers - $D_{density}$ and D_{color} and training steps. These parameters are provided in Table 1 for each model. In the small, base light and base versions, TriNeRFLet uses the exact same MLPs as in INGP. TriNeRFLet is trained using the Adam optimizer with a learning rate of 0.01 and an exponential decay scheduler. The wavelet low frequencies in LL (of size $N_{LL} \times N_{LL}$) are randomly initialized, and the other frequencies in LH, HL and HH at all levels are initialized with zeros. The other training details are very similar to INGP.

Since the wavelet reconstruction component affects the runtime of TriNeRFLet, we start the training with lower resolutions of the planes, training coarse to fine. The initial plane resolution is set to N_{base} and it is increased throughout the training till it reaches the size N_{final} that has L wavelet levels. As described above, due to the wavelet multiscale structure, when we move to the next resolution in training we simply add another level (of higher frequencies) to the wavelet representation and add it to the training, where the lower levels are kept as is (but continue to train).

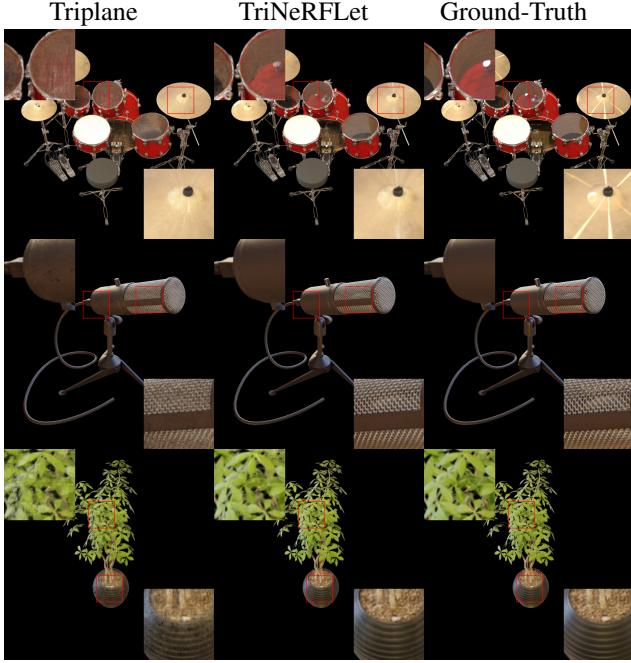


Figure 7. *NeRF reconstruction qualitative results*. Notice the improvement in reconstruction quality of TriNeRFLet compared to Triplane. More visual results appear in the project page.

TriNeRFLet implementation is built upon [39], which is a Pytorch-based implementation of INGP.

To evaluate the performance of our method, we use the Blender dataset from [27], which contains 8 synthetic 3D scenes with complex geometry and lighting. We compare TriNeRFLet with the vanilla NeRF [27], regular Triplane (with base config), INGP [28] and 3D Gaussian Splatting [20]. Results are reported in Table 2, and qualitative examples are presented in Figure 7. We ran all TriNeRFLet versions on a single A6000 GPU. The reported results indeed demonstrate the improvement in performance that the multiscale wavelet brings, as it outperforms regular Triplane by a significant margin, and improves over the current SOTA.

Regarding rendering time, the TriNeRFLet (all versions) frames per second (FPS) are approximately 1 – 4, while INGP and 3D Gaussian Splatting are approximately 0.6 – 1.2 and 135 respectively. These FPS are tested on a single A6000 GPU. Note that TriNeRFLet rendering FPS is compatible and even slightly better than INGP, which is known to be fast within the NeRF approaches. 3D Gaussian Splatting, which is not a NeRF method, is faster in its rendering time since it uses a different internal structure than NeRF. When compared only to NeRF-based solutions, our method is very fast with competitive reconstruction performance. Note that although our rendering time is faster, the training time of TriNeRFLet is higher than INGP. This is due to wavelet reconstruction, which needs to be done at each

Method	100 → 400		200 → 800	
	PSNR↑	LPIPS↓	PSNR↑	LPIPS↓
NeRF*	25.56	0.170	27.47	0.128
NeRF-Bi	24.74	0.244	26.67	0.175
NeRF-Liif	25.36	0.125	27.34	0.096
NeRF-Swin	24.85	0.108	26.78	0.086
NeRF-SR (SS)	<u>28.07</u>	<u>0.071</u>	28.46	0.076
TriNeRFLet-SR	28.55	0.061	<u>29.49</u>	<u>0.051</u>
NVSR (3D supervised)	**		29.53	0.045

Table 3. *Blender ×4 Super Resolution Results* on low resolution input of 100×100 and 200×200 . NVSR is trained with multiview (3D) supervision and therefore we consider it as upper-bound. **Bold** is best, underline is second. **For $100 \rightarrow 400$ NVSR reports results of only 4 shapes in the Blender dataset. In this case, we outperform it with PSNR of 29.18dB for our method compared to 28.5dB of NVSR.

training step, while only once during rendering time.

5.2. Blender Scene Super-Resolution

We turn to present our results for TriNeRFLet super-resolution. The objective here is to generate novel views with $\times 4$ resolution as input images. To this end, we tested two different settings on the Blender dataset [27], 100 to 400 and 200 to 800. For the first one, we use TriNeRFLet with $N_{LR} = 256$ and $N = 1024$, while in the second setting $N_{LR} = 512$ and $N = 2048$. In both settings $Ch = 16$. We use [12] as a base for implementing TriNeRFLet super-resolution, as it is more suitable for this task.

We compare our method with NeRF*, NeRF-Bi, NeRF-Liif, NeRF-Swin, NeRF-SR (SS) (these methods do not require 3D supervision and are described in [40]) and NVSR (results of $200 \rightarrow 800$ were provided to us by the authors) [4]. Results are reported in Table 3. Note that NVSR requires 3D supervision (paired LR and HR multiview images) and therefore it is considered only as a reference. When compared only to methods that use 2D supervised methods, we perform the best in most of the cases. Qualitative examples of our approach appear in Figure 8. More super-resolution results appear in the sup. mat, including results on LLFF dataset.

5.3. LLFF Scene Super-Resolution

We turn to demonstrate the TriNeRFLet super-resolution scheme on the LLFF dataset [26], which contains 8 real-world captured scenes. As in the Blender super-resolution case, we compare our method with the same methods but for input resolution 378×504 and $\times 4$ upscaling. We use the same settings as in the Blender $200 \rightarrow 800$ experiment except for the guidance scale, which we choose to be 2.5 instead of 7.5. Following [40], we trained TriNeRFLet using the low-resolution version of all images and then tested the super-resolution results also on all the images. Simi-

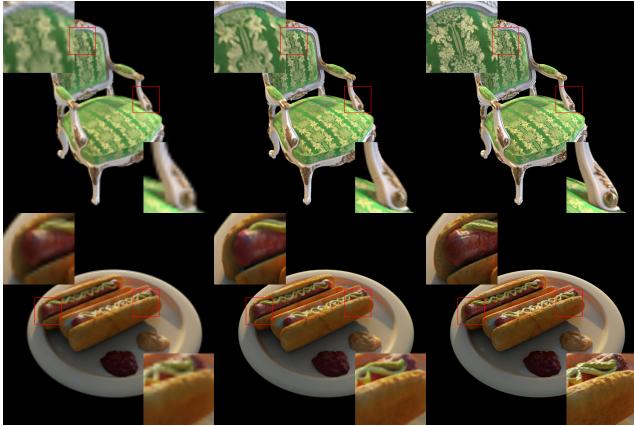


Figure 8. *TriNeRFLet SR qualitative results.* (Left) LR image; (Middle) TriNeRFLet SR; (Right) Ground truth HR image.

Method	PSNR↑	LPIPS↓
NeRF*	24.47	0.388
NeRF-Bi	23.90	0.481
NeRF-Liif	24.76	0.292
NeRF-Swin	23.26	0.247
NeRF-SR (SS)	25.13	<u>0.244</u>
TriNeRFLet-SR	<u>25.00</u>	0.203

Table 4. *LLFF × 4 Super Resolution Results* for low resolution input of 378×504 . **Bold** is best, underline is second.

lar to the Blender experiment, we ran LLFF TriNeRFLet super-resolution on a single A6000 RTX GPU, where each scene runs almost for 12 hours. Table 4 reports the results. Our method is almost as good as [40] in PSNR and outperforms it in the LPIPS metric by a margin. Our LPIPS performance indicates that our method manages to create a high-resolution scene with better details and that is perceptually more similar to the high-resolution scene.

6. Conclusions

This paper introduced a new NeRF structure that relies on multiscale wavelet representation. This structure improved the performance of Triplane, which lagged behind NeRF SOTA methods, achieving competitive results. Having the multiscale structure enabled us to implement a new diffusion-based super-resolution for NeRF.

Despite the advantages our proposed approach brings, it also has some limitations, mainly due to the additional runtime overhead it brings to training time. This can be partially mitigated by using its light-weight versions and using the coarse to fine training, which reduces training time.

Our proposed approach of learning the features in the wavelet domain instead of the original one, opens the door for many future works to implement it in other applications.

We believe this approach can benefit tasks where data is sparse or training suffers from over-fitting, by regularizing the high frequencies in the wavelet domain.

References

- [1] Shady Abu-Hussein, Tom Tirer, and Raja Giryes. Adir: Adaptive diffusion for image reconstruction. *arXiv preprint arXiv:2212.03221*, 2022. 4
- [2] Omri Avrahami, Ohad Fried, and Dani Lischinski. Blended latent diffusion. *arXiv preprint arXiv:2206.02779*, 2022. 4
- [3] Omri Avrahami, Dani Lischinski, and Ohad Fried. Blended diffusion for text-driven editing of natural images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18208–18218, 2022. 4
- [4] Yuval Bahat, Yuxuan Zhang, Hendrik Sommerhoff, Andreas Kolb, and Felix Heide. Neural volume super-resolution. *arXiv preprint arXiv:2212.04666*, 2022. 2, 3, 4, 6, 8
- [5] 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. 1
- [6] 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
- [7] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16123–16133, 2022. 1, 3
- [8] Hyungjin Chung, Eun Sun Lee, and Jong Chul Ye. Mr image denoising and super-resolution using regularized reverse diffusion. *IEEE Transactions on Medical Imaging*, 42(4):922–934, 2023. 4
- [9] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021. 4
- [10] Rinon Gal, Dana Cohen Hochberg, Amit Bermano, and Daniel Cohen-Or. Swagan: A style-based wavelet-driven generative model. *ACM Trans. Graph.*, 40(4), 2021. 4
- [11] Tiantian Guo, Tongpo Zhang, Enggee Lim, Miguel Lopez-Benitez, Fei Ma, and Limin Yu. A review of wavelet analysis and its applications: Challenges and opportunities. *IEEE Access*, 10:58869–58903, 2022. 2, 5
- [12] Yuan-Chen Guo, Ying-Tian Liu, Ruizhi Shao, Christian Laforte, Vikram Voleti, Guan Luo, Chia-Hao Chen, Zi-Xin Zou, Chen Wang, Yan-Pei Cao, and Song-Hai Zhang. threestudio: A unified framework for 3d content generation. <https://github.com/threestudio-project/threestudio>, 2023. 8
- [13] Florentin Guth, Simon Coste, Valentin De Bortoli, and Stephane Mallat. Wavelet score-based generative modeling. In *Advances in Neural Information Processing Systems*, pages 478–491, 2022. 4

- [14] Ayaan Haque, Matthew Tancik, Alexei Efros, Aleksander Holynski, and Angjoo Kanazawa. Instruct-nerf2nerf: Editing 3d scenes with instructions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023. 6
- [15] Amir Hertz, Or Perel, Raja Giryes, Olga Sorkine-Hornung, and Daniel Cohen-Or. Sape: Spatially-adaptive progressive encoding for neural optimization. *Advances in Neural Information Processing Systems*, 34:8820–8832, 2021. 4
- [16] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 4
- [17] Zahra Kadkhodaie, Florentin Guth, Stéphane Mallat, and Eero P Simoncelli. Learning multi-scale local conditional probability models of images. In *The Eleventh International Conference on Learning Representations*, 2023. 4
- [18] Bahjat Kawar, Michael Elad, Stefano Ermon, and Jiaming Song. Denoising diffusion restoration models. In *Advances in Neural Information Processing Systems*, 2022. 4
- [19] Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. Imagic: Text-based real image editing with diffusion models. *arXiv preprint arXiv:2210.09276*, 2022. 4
- [20] 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. 2, 3, 8
- [21] Ming Li, Pan Zhou, Jia-Wei Liu, Jussi Keppo, Min Lin, Shuicheng Yan, and Xiangyu Xu. Instant3d: Instant text-to-3d generation, 2023. 2, 4
- [22] Chen-Hsuan Lin, Wei-Chiu Ma, Antonio Torralba, and Simon Lucey. Barf: Bundle-adjusting neural radiance fields. *arXiv preprint arXiv:2104.06405*, 2021. 3
- [23] David B. Lindell, Dave Van Veen, Jeong Joon Park, and Gordon Wetzstein. Bacon: Band-limited coordinate networks for multiscale scene representation. In *CVPR*, 2022. 4
- [24] Stéphane Mallat. *A Wavelet Tour of Signal Processing, Third Edition: The Sparse Way*. Academic Press, Inc., 2008. 4, 11
- [25] 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 *CVPR*, 2021. 2
- [26] Ben Mildenhall, Pratul P Srinivasan, Rodrigo Ortiz-Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi, Ren Ng, and Abhishek Kar. Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Transactions on Graphics (TOG)*, 38(4):1–14, 2019. 8
- [27] 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, 8
- [28] 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. 1, 2, 3, 5, 8
- [29] Keunhong Park, Utkarsh Sinha, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Steven M. Seitz, and Ricardo Martin-Brualla. Deformable neural radiance fields. *arXiv preprint arXiv:2011.12948*, 2020. 3
- [30] Hao Phung, Quan Dao, and Anh Tran. Wavelet diffusion models are fast and scalable image generators. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10199–10208, 2023. 4
- [31] Qiang Qiu, Xiuyuan Cheng, Guillermo Sapiro, et al. Dcfnet: Deep neural network with decomposed convolutional filters. In *International Conference on Machine Learning*, pages 4198–4207, 2018. 4
- [32] Daniel Rho, Byeonghyeon Lee, Seungtae Nam, Joo Chan Lee, Jong Hwan Ko, and Eunbyung Park. Masked wavelet representation for compact neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20680–20690, 2023. 4
- [33] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 2, 4, 5, 6
- [34] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35:36479–36494, 2022. 2
- [35] Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad Norouzi. Image super-resolution via iterative refinement. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022. 4
- [36] Vishwanath Saragadam, Daniel LeJeune, Jasper Tan, Guha Balakrishnan, Ashok Veeraraghavan, and Richard G Baraniuk. Wire: Wavelet implicit neural representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18507–18516, 2023. 4
- [37] J Ryan Shue, Eric Ryan Chan, Ryan Po, Zachary Ankner, Jiajun Wu, and Gordon Wetzstein. 3d neural field generation using triplane diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20875–20886, 2023. 2, 4
- [38] 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. *NeurIPS*, 2020. 2
- [39] Jiaxiang Tang. Torch-npg: a pytorch implementation of instant-npg, 2022. <https://github.com/ashawkey/torch-npg>. 8
- [40] 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 30th ACM International Conference on Multimedia*, pages 6445–6454, 2022. 8, 9
- [41] Jay Whang, Mauricio Delbracio, Hossein Talebi, Chitwan Saharia, Alexandros G Dimakis, and Peyman Milanfar. Deblurring via stochastic refinement. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16293–16303, 2022. 4

- [42] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelNeRF: Neural radiance fields from one or few images. In *CVPR*, 2021. 2
- [43] Jason J. Yu, Konstantinos G. Derpanis, and Marcus A. Brubaker. Wavelet Flow: Fast training of high resolution normalizing flows. In *NeurIPS*, 2020. 4

A. Ablation

For all the experiments we conducted in the paper, we applied the Biorthogonal 6.8 (Bior6.8) wavelet. We present here an ablation experiment that checks the impact of other wavelet filters. We compare the reconstruction performance of 4 scenes from the Blender dataset. The setups we tested are vanilla Triplane (no wavelet), Haar, Biorthogonal 2.2 (Bior2.2), Biorthogonal 2.6 (Bior2.6), Biorthogonal 4.4 (Bior4.4) and Biorthogonal 6.8 (Bior6.8). The results are presented in Figure 9. We use the TriNeRFLet Base Light setting. Vanilla Triplane has the same size and structure as this setting. Note that all wavelets achieve a significant performance advantage over vanilla Triplane (as shown also in the paper). Generally, higher order wavelets provide a better representation of smooth functions [24]. Thus, it is not surprising that Bior6.8, Bior4.4 and Bior 2.6 achieve better performance than Bior 2.2 and Haar. Note though that in terms of training time, training with Haar wavelet is faster by up to 30% compared to the other wavelets as Haar complexity is $O(N)$ and the complexity of applying the other wavelets is $O(N \log(N))$.

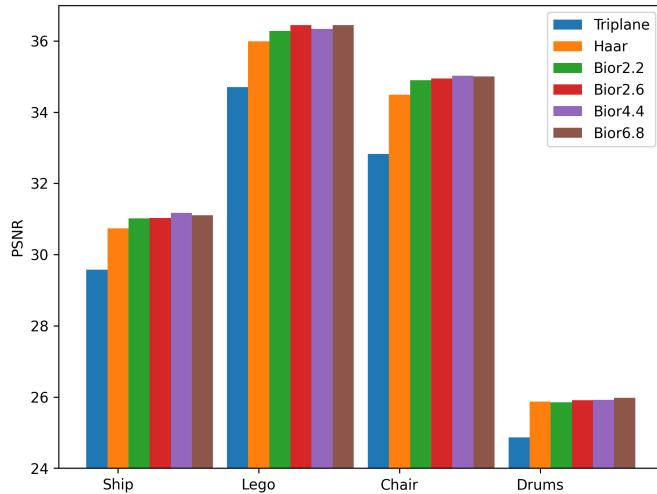


Figure 9. Performance comparison between different wavelet filters and no wavelet at all (Triplane).