

PlückerRF: A Line-based 3D Representation for Few-view Reconstruction

Sam Bahrami
 The Australian National University
 Canberra, Australia
 sam.bahrami@anu.edu.au

Dylan Campbell
 The Australian National University
 Canberra, Australia
 dylan.campbell@anu.edu.au

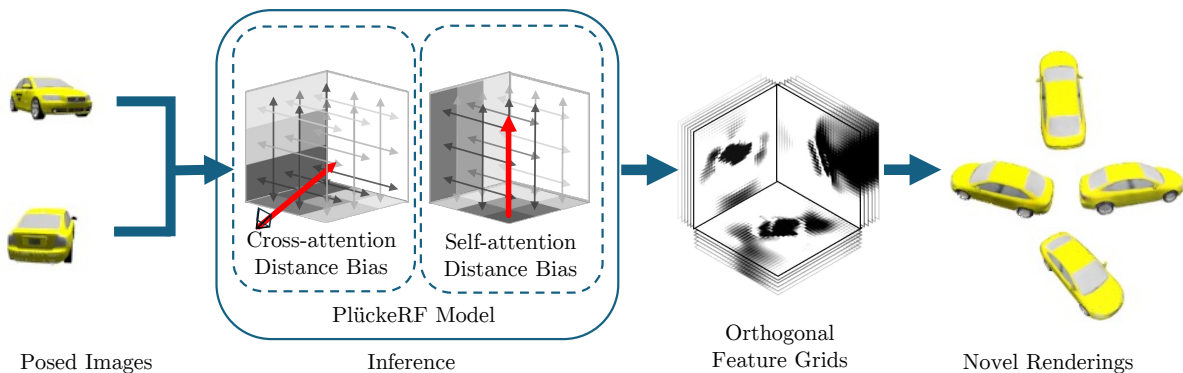


Figure 1. PlückerRF: A feed-forward method that directly predicts a 3D representation from posed input images. The model uses a line–line distance biased attention mechanism that helps the model interpret which parts of the images are important for constructing specific areas of the 3D output. The resulting 3D representation is structured as three orthogonal feature grids, enabling efficient novel view synthesis.

Abstract

Feed-forward 3D reconstruction methods aim to predict the 3D structure of a scene directly from input images, providing a faster alternative to per-scene optimization approaches. Significant progress has been made in single-view and few-view reconstruction using learned priors that infer object shape and appearance, even for unobserved regions. However, there is substantial potential to enhance these methods by better leveraging information from multiple views when available. To address this, we propose a few-view reconstruction model that more effectively harnesses multi-view information. Our approach introduces a simple mechanism that connects the 3D representation with pixel rays from the input views, allowing for preferential sharing of information between nearby 3D locations and between 3D locations and nearby pixel rays. We achieve this by defining the 3D representation as a set of structured, feature-augmented lines—the PlückerRF representation. Using this representation, we demonstrate improvements in reconstruction quality over the equivalent triplane representation and state-of-the-art feedforward reconstruction methods.

1. Introduction

3D reconstruction is the task of estimating the shape and appearance of objects or scenes from data. This typically involves creating a 3D model from a set of 2D images, which may include additional information such as camera pose or depth maps. With enough images, methods like Neural Radiance Fields (NeRF) [20] and more recently Gaussian Splatting [12] allows one to reconstruct a 3D scene through novel view synthesis (NVS). The challenge comes when there is limited data, which is the case for few-view and single-view 3D reconstruction tasks.

Initial neural representations for NVS were based on fully-connected multi-layer perceptron (MLP) models [20]. Since then interpretable neural representations have been developed, such as voxel feature grids [16], hashing [21], hybrid implicit–explicit triplane representations [1, 3], and Gaussian mixtures [12]. Each of these representations presents its own trade-offs in terms of computational efficiency, memory usage, and rendering quality.

Voxel grids and triplane representations are learnable 3D structures that encode information in the form of feature grids and orthogonal feature planes, respectively. To infer a novel view, a point in 3D space is queried through this

representation to obtain its corresponding feature encoding, which is processed by a small neural network to determine the pixel’s color and density. Unlike MLP representations, voxel grids and triplanes offer an inherently interpretable structure, enabling the application of useful inductive biases, such as convolutions that share information between neighboring points. This structured nature allows voxel grids and triplanes to be inferred through image-to-image methods, such as encoder–decoder style feed-forward neural networks [8, 29].

Several single-view 3D reconstruction methods have been adapted to work with additional, but still few, views [31, 41]. This is often referred to as few-view reconstruction. Often this adaptation is done *without* providing additional context about the relationships between the views. For example, by inferring an explicit representation from each individual view and merging the individual outputs [42], simply concatenating additional image and pose encodings to a model’s input [11, 15], or averaging feature representations across individual image and pose encodings [41].

Previous work has shown that incorporating an image pixel ray-to-ray distance bias can improve the quality of few-view 3D reconstructions [32]. We define a distance measure that quantifies the proximity between a camera ray for an image pixel, and a ray orthogonal to a pixel of our orthogonal feature grid representation—the PlückerRF representation. Explicitly we define this as a line-to-line distance measure, by converting these rays to origin invariant line representations. We use these distances to bias the attention mechanism in a feed-forward transformer model, encouraging information sharing between lines that are close or intersecting while penalizing sharing between distant lines. The intuition for this is that image regions corresponding to specific areas of the 3D representation should exert a stronger influence on those areas within the model. A high level summary of our approach can be seen in Fig. 1.

By incorporating additional contextual information from a few supplementary views (two views in our experiments), our approach enhances the accuracy and realism of the novel views inferred from our 3D neural field. Our contributions are

1. the PlückerRF representation that facilitates geometrically-meaningful sharing of information within the 3D representation and between it and the available input images; and
2. a simple few-view reconstruction model that uses this internal representation to predict a 3D neural field in a single forward pass.

2. Related Work

Commonly, 3D reconstructions are represented using point clouds [38], meshes [34], neural representations [3, 20] or

Gaussian splats [12]. Models designed for single or few-view reconstruction typically learn a prior from a dataset, and with sufficiently large datasets, this prior can enable robust generalization.

Neural 3D Reconstruction and Representations. A neural radiance field (NeRF) [20] is a machine learning model that learns the mapping between a 3D point and view direction with a corresponding color and density. The original NeRF approach employs a multi-layer perceptron (MLP) neural network to define an implicit function which takes a point and a ray direction, and returns a color and density. To render a posed image, a ray is cast through each pixel. Multiple positions along each ray are then queried using the model, informed by the camera’s pose, and volumetric rendering is applied to compute the final pixel color. Hashing methods have addressed some of the drawbacks of the original NeRF, by using a multi-resolution hash table encoding that efficiently represents spatial data and leveraging hierarchical representations to quickly capture both global and fine-grained details in the scene [21].

Subsequent works implemented alternative structured representations such as pixel-aligned Gaussian mixtures [2, 28, 31] based on Gaussian splats [12], voxel grids [16, 29], and triplane feature grids [1, 3, 6]. These structured representations offer memory-efficient 3D representations of the scene, at a trade-off with the representation resolution. We use a triplane-like representation, as it provides greater resolution with a smaller memory footprint compared to a voxel representation, and we can make use of the structured representation to guide our PlückerRF model.

Single-View 3D Object Reconstruction. 3D representations such as point clouds [37, 38], meshes [39], signed distance fields [24], NeRFs [8, 9, 22, 41], and Gaussian Splats [31] have dominated learning-based single-view reconstruction approaches. These methods learn a prior representation across a dataset, which generalizes to new objects or scenes [1, 22, 30–32]. Initial single-view NeRF methods parameterize the MLP through latent representations from the source image [9, 41]. Later works used structured NeRF representations such as voxel feature grids [29], and hybrid implicit–explicit triplane representations [1, 8]. With triplane representations, large feed-forward reconstruction models have become a popular method for single and few-view 3D reconstruction [8]. These models use large scale datasets (e.g., 800k+ objects [5]) which helps generalize across categories.

Another approach has been to use generative models to convert the single–view reconstruction to a multi-view reconstruction problem. These approaches may fine-tune their 3D representation using generated multi–view instances [17] or use estimated images from other views for

guidance in image-to-3D reconstruction [6, 19].

Most recently pixel aligned Gaussian mixture representations have also shown strong performance in single-view reconstruction, taking advantage of the very fast rendering to speed up training, requiring less compute resources to get results comparable to those using the triplane representations [30, 31].

Our PlückerRF approach is based on a feed-forward reconstruction model [8, 31]. These models are very simple, often comprised of purely transformer layers, and by using an orthogonal feature grid representation, like a triplane, our model can infer the complete feature representation in a forward pass.

Few-View 3D Object Reconstruction. Few-view often means a number of input views between two [31, 41] and six [32]. Utilizing multiple views provides more context to the model in two ways, first by reducing the amount of unobserved object surface area, and second by making it possible for a model to triangulate 3D surface geometry. Many of the previous single view reconstruction tasks have been adapted to use multiple views. Some of these methods simply merge representations from each individual view [31, 41, 42], or provide the additional image and camera pose to their model as concatenated image encoding tokens [15].

Other approaches leverage geometric relationships between images to enhance 3D reconstruction. Even without known camera poses, it is possible to predict the relative 3D poses of unposed 2D images and estimate the scaled depth of associated points in each image [14, 33, 35, 36]. This serves as a valuable starting point for reconstructing scenes from image pairs [28]. Another approach jointly optimizes object shape and camera poses from few-view inputs, reducing reliance on precise pose initialization and yielding consistent 3D reconstructions [10, 40].

When camera poses are known, this information can be utilized to determine which parts of the reconstruction each individual image should focus on [32]. 3D representation-free novel view synthesis models have emerged that directly map sparse views to high-quality novel images without explicit 3D representations [11]. Many of these multi-view reconstruction methods use Plücker coordinates as a positional encoding of their inputs [4, 11, 27].

In contrast, our work utilizes the pose information from the few views, and geometric information from our structured feature representation, using a Plücker coordinate representation for both. Pixels in our input images can be attended to based on their relationship with the underlying 3D representation—our PlückerRF representation.

3. PlückerRF Few-view Reconstruction

The goal of this work is to render novel viewpoints of previously unseen objects from a few posed images. To do so,

we train a transformer-based neural network to infer three orthogonal feature grids from the provided input views. We begin by providing background information on Plücker coordinates and the feature grid representation. Next, we provide a high level overview of our model, followed by details on our 2D image representation and our geometric feature grid representation, explaining how they interact within our transformer attention mechanism. We conclude with detail on our pretrained image encoder and training information.

3.1. Preliminaries

Orthogonal Feature Grid Representation. We use three orthogonal, axis-aligned feature grids that are aligned with the world coordinate system to model the opacity and radiance features at every 3D point. This is also known as a triplane representation [1]. Concretely, this representation is defined as a set of three orthogonal axis-aligned feature grids $\mathcal{T} = \{\mathbf{T}_{xy}, \mathbf{T}_{yz}, \mathbf{T}_{zx}\}$ where $\mathbf{T}_{xy}, \mathbf{T}_{yz}, \mathbf{T}_{zx} \in \mathbb{R}^{M \times M \times d_T}$ with M being the square grid width and height respectively, and d_T is the feature dimension. The grid covers the dilated object bounding box $[-1, 1]^3$, which defines the 3D space the object exists in. For details on how this representation is used to render novel views, see Sec. 3.7.

Plücker Coordinates. Plücker coordinates are a line representation $\mathbf{l} = (\mathbf{d}, \mathbf{m}) \in \mathbb{P}^5$ (projective 5-space) that may be computed from a ray with origin $\mathbf{o} \in \mathbb{R}^3$ and direction $\mathbf{d} \in \mathbb{S}^2$ by taking the cross product [25]:

$$\mathbf{l} = (\mathbf{d}, \mathbf{m}) = (\mathbf{d}, \mathbf{o} \times \mathbf{d}) \quad (1)$$

They are invariant to the choice of origin along the ray and are homogeneous coordinates for the line, that is, $(\lambda \mathbf{d}, \lambda \mathbf{m})$ for $\lambda \neq 0$ represents the same line. We fix $\lambda = 1$. Given two lines represented by Plücker coordinates $\mathbf{l}_1 = (\mathbf{d}_1, \mathbf{m}_1)$ and $\mathbf{l}_2 = (\mathbf{d}_2, \mathbf{m}_2)$, the closest distance between the lines is given by

$$d(\mathbf{l}_1, \mathbf{l}_2) = \begin{cases} |\mathbf{d}_1^\top \mathbf{m}_2 + \mathbf{d}_2^\top \mathbf{m}_1| / \|\mathbf{d}_1 \times \mathbf{d}_2\|_2 & \text{if } \mathbf{d}_1 \times \mathbf{d}_2 \neq 0 \\ \|\mathbf{d}_1 \times (\mathbf{m}_1 - (\mathbf{d}_1^\top \mathbf{d}_2) \mathbf{m}_2)\|_2 & \text{otherwise.} \end{cases} \quad (2)$$

When $\mathbf{d}_1 \times \mathbf{d}_2 \neq 0$, the lines are not parallel, allowing us to compute the shortest distance between skew lines. In the other case, the lines are parallel and the distance is calculated by projecting the moment vectors onto a plane perpendicular to their common direction.

3.2. Images-to-Feature-Grid Transformer Decoder

A visual summary of our method can be seen in Fig. 2. Our method infers three orthogonal feature grids from the provided input views. We use a transformer decoder architecture with cross-attention and self-attention layers within each transformer block. The initial input to our model is a

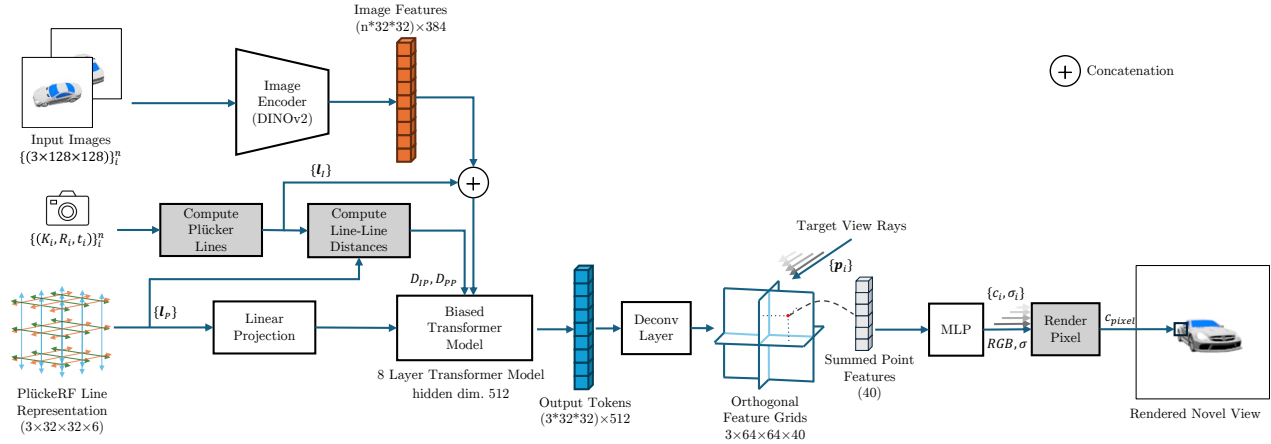


Figure 2. Flowchart for the PlückerRF model. We use a pre-trained vision transformer (DINOv2 [23]) to encode the input images (Sec. 3.6). We take the camera information from the images and create rays corresponding with the image patches from our image encoder, which we convert to Plücker coordinates (Sec. 3.3). Our final input is the set of lines corresponding to the pixels in the orthogonal feature grid representation, which we call the PlückerRF line representation (Sec. 3.4). We calculate the distance between these Plücker coordinates and those from our input images, and apply these as an attention bias in our transformer attention blocks (Sec. 3.5). The output of the transformer is reshaped to three orthogonal feature planes, which we use for novel view synthesis (Sec. 3.7). Our network is trained end-to-end with image reconstruction losses (Sec. 3.8).

linear projection of the PlückerRF line parameters (Sec. 3.4) to the transformer hidden dimension d_D . This corresponds with the shape of the three feature grid representation, which are internally represented as flattened feature grids of $3N^2 \times d_D$ dimension.

Our few-view images are applied to the transformer model through the cross-attention layers, where each image is pre-processed into patch-wise image features from a pretrained image encoder, detailed in Sec. 3.6. Plücker coordinates for each patch are computed from the camera information and concatenated to each pixel patch in its feature dimension, as done in other works as a positional encoding [4, 11, 27]. In addition to patch-wise image features, the DINOv2 encoder provides an additional CLS token. Since the CLS token does not correspond with any particular part of the image, we set the distance between it and all of the PlückerRF lines to 0—the CLS token should incur no distance-based attention penalty. The set of n input images are then concatenated together.

The resulting output of our transformer model is reshaped to our three orthogonal feature grids. We process this through a deconvolution layer to upsample the representation, i.e., from $(3 \times N \times N \times d_D)$ to $(3 \times M \times M \times d_T)$, where d_T is our orthogonal grids feature dimension, and $M = 2N$. These feature grids are then used to render the input views and unseen views of our object, using standard NeRF volumetric rendering. These views are compared with ground-truth images to calculate a training loss and optimize the model architecture. The

model hence learns a prior for the training dataset.

3.3. Image Plücker Coordinates

Our model processes a set of n input images $\{I_i\}_{i=1}^n$. These images are pre-processed into feature vectors of pixel patches using a pretrained DINOv2 [23] image encoder (see Sec. 3.6). We utilize Plücker coordinates to represent lines from the camera center passing through the middle of the pixel patches in the encoded image. These rays are generated using the camera intrinsics $\mathbf{K} \in \mathbb{R}^{3 \times 3}$ and extrinsics $\mathbf{R} \in \mathbb{R}^{3 \times 3}$, $\mathbf{t} \in \mathbb{R}^3$ relative to the first input camera for each image in our set.

Let $\mathbf{u} = (u, v, 1)^T$ be the homogeneous coordinates of a pixel center in the image. The corresponding ray direction \mathbf{d} in the world coordinate system is given by $\mathbf{d} = (\mathbf{R}\mathbf{K}^{-1}\mathbf{u}) / \|\mathbf{K}^{-1}\mathbf{u}\|$. The origin of the ray \mathbf{o} is the camera center, given by the translation vector $\mathbf{o} = \mathbf{t}$. We convert these rays $\mathbf{r} = (\mathbf{o}, \mathbf{d})$ to Plücker coordinates \mathbf{l} following Eq. (1). This results in a set of Plücker coordinates corresponding to each pixel patch of each input image, denoted as $\{I_i\}_{i=1}^n$ for the lines from each image i in our image set.

The conversion to Plücker coordinates abstracts away the specific location of the camera center, focusing instead on the direction of the rays. Plücker coordinates represent lines based solely on their direction and moment, independent of any specific point along the line, thereby ensuring that translations along the ray do not change the encoding. This invariance to the camera’s position along the line is advantageous in our model, as moving the camera forward or backward along a pixel ray should not significantly affect the

color or geometry associated with that particular pixel. As a result, the representation is robust to variations in the camera’s position along the view direction. The invariance also simplifies the model’s learning process by eliminating the need to identify different inputs (camera positions) as being the same (defining the same ray).

3.4. The PlückerRF Representation

Our model infers three orthogonal feature grids which are used for neural rendering. We define rays orthogonal to each grid plane through each pixel center of our orthogonal feature grids, at the resolution of our model’s internal representation, i.e., $3 \times N \times N$ rays. We convert these rays to lines following Eq. (1), and denote this set of lines as $\{\mathbf{l}_{P_i}\}_{i=1}^{3N^2}$ —our PlückerRF line representation. We use these to establish a direct geometric relationship between the feature grid features and the input image features through a line–line distance-biased attention. The input tokens to our biased transformer are instantiated as a linear projection of $\{\mathbf{l}_{P_i}\}_{i=1}^{3N^2}$ to the model hidden dimension d_D .

3.5. Line–line Distance-biased Attention

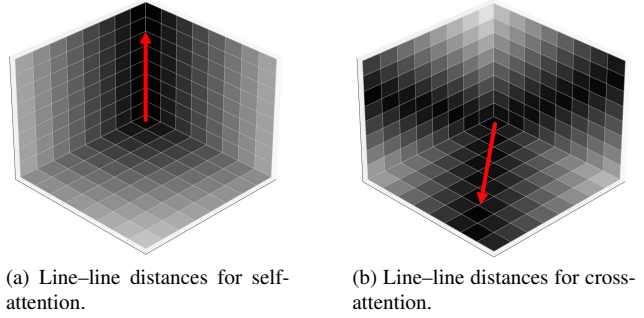
We add a distance bias to the attention mechanism, based on the line-to-line distance between each corresponding image patch $\{\mathbf{l}_{I_i}\}_{i=1}^n$ for our n images, and the lines of our PlückerRF representation $\{\mathbf{l}_{P_i}\}_{i=1}^{3N^2}$. The intuition is that (a) distant lines should have low attention weights, and hence have a high distance penalty, (b) nearby lines should have a smaller penalty, and hence higher attention weights, and (c) intersecting lines should not be penalized at all (see Fig. 3). This encourages information to be shared between neighboring and intersecting PlückerRF pixels and reduces sharing between unrelated parts of the representation. We observe that the latter might be sub-optimal in the presence of symmetries and so provide a mechanism to attenuate this bias if it is unhelpful (see γ , below).

The standard scaled dot product attention mechanism in a transformer cross-attention layer is calculated as follows, using the query \mathbf{Q} , key \mathbf{K} , and value \mathbf{V} matrices,

$$\text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_I}} \right) \mathbf{V}, \quad (3)$$

where the keys and values $\mathbf{K}, \mathbf{V} \in \mathbb{R}^{nE^2 \times d_I}$ come from our encoded image features, for n views, image encodings of resolution E , and feature dimension d_I . The query matrix $\mathbf{Q} \in \mathbb{R}^{3N^2 \times d_I}$ originates from a linear projection of the PlückerRF line representation—the 3D representation derived from our orthogonal feature grids.

In our model, we bias the attention matrix with the distance between our lines representing the keys $\mathbf{l}_k \in \{\mathbf{l}_{I_i}\}_{i=1}^n$ and our line-based 3D representation representing



(a) Line–line distances for self-attention.

(b) Line–line distances for cross-attention.

Figure 3. A visual representation of the Plücker coordinate distance between a line from our transformer keys (represented by the red arrow), and the lines corresponding to the queries—the PlückerRF line representation. The shading on the squares on each axis-aligned plane represents the distance between $\{\mathbf{l}_{P_i}\}_{i=1}^{3N^2}$ and the red arrow. In practice, those with a greater distance (lighter squares) will be penalized in attending to this feature. (a) The self-attention case, where the red arrow is one of the lines $\{\mathbf{l}_{P_i}\}_{i=1}^{3N^2}$ from the 3D PlückerRF representation. (b) The cross-attention case, where the red arrow corresponds to a ray from an input camera passing through an image pixel patch.

the queries $\mathbf{l}_q \in \{\mathbf{l}_{P_i}\}_{i=1}^{3N^2}$,

$$\text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_I}} - \gamma d(\mathbf{l}_q, \mathbf{l}_k) \right) \mathbf{V}, \quad (4)$$

where the distance is calculated using Equation (2), and $\gamma \in \mathbb{R}_+$ is a learnable parameter controlling the relative importance of the distance bias. This provides the model an inductive bias to attend to image features which intersect with areas of the feature grid representation, see Fig. 3b.

In the self-attention layers, we compare each line in $\{\mathbf{l}_{P_i}\}_{i=1}^{3N^2}$ with all other lines in the set for its distance calculation, for every pixel in our feature grids. This provides an inductive bias similar to that of a convolutional layer in a neural network, where neighboring parts of the image share information, but extends this to a 3D locality bias *while only using 2D data structures* (see Fig. 3a). This contrasts with 3D convolutional layers, which have a locality bias, but require processing a volumetric data structure, which is less memory and compute efficient. Unlike convolutional layers, our approach retains the flexibility of transformers, which can attend across the entire feature map to learn global and local relationships. Additionally, the learnable parameter γ allows the bias to be ignored if it is not useful.

We denote these pairwise distance matrices as \mathbf{D}_{IP} for the cross-attention and \mathbf{D}_{PP} for the self-attention cases.

3.6. Image Encoder

Given RGB images as an input, our method first applies a pre-trained vision transformer (ViT) [13] to encode the

image to patch-wise feature tokens. We initialize this as a pretrained DINOv2 [23] model. The model uses 14×14 pixel image patches.

3.7. Neural Rendering

We follow the triplane NeRF formulation which was introduced by Chan et al. [1]. We use an MLP to predict a color and density from the point features queried from the orthogonal feature grid representation $\mathcal{T} = \{\mathbf{T}_{xy}, \mathbf{T}_{yz}, \mathbf{T}_{zx}\}$.

We compute rays corresponding to each target pixel for our target view given its camera information. We sample r points $\{\mathbf{p}_i\}_{i=1}^r$ along the ray to determine the feature representation, and repeat this ray marching process for every pixel in our rendered view.

For neural rendering, each point $\mathbf{p} \in \mathbb{R}^3$ is projected onto each plane to obtain three feature vectors, $\mathbf{f}_{xy}, \mathbf{f}_{yz}, \mathbf{f}_{zx}$. For example, feature vector $\mathbf{f}_{xy} \in \mathbb{R}^{d_T}$ is obtained by projecting \mathbf{p} onto the plane and bilinearly interpolating the grid feature of \mathbf{T}_{xy} at that location. That is $\mathbf{f}_{xy} = \text{interp}(\text{proj}_{xy}(\mathbf{p}), \mathbf{T}_{xy})$ for $\text{proj}_{xy}(\mathbf{p}) = [p_x, p_y]$ and interp being the bilinear interpolation function. This is similarly done to obtain $\mathbf{f}_{yz}, \mathbf{f}_{zx}$.

These features are summed together and processed by a small decoder MLP network to interpret the features as color (RGB) and density σ . Our features $\mathbf{f}_{\mathbf{p}}$ at point $\mathbf{p} = (p_x, p_y, p_z)$ are defined by $\mathbf{f}_{\mathbf{p}} = \mathbf{f}_{xy} + \mathbf{f}_{yz} + \mathbf{f}_{zx}$, where $\mathbf{f}_{\mathbf{p}} \in \mathbb{R}^{d_T}$. These are processed by an MLP to predict the color and density at that point,

$$(c_i, \sigma_i) = f(\mathbf{p}_i) = \text{MLP}(\mathbf{f}_{\mathbf{p}}). \quad (5)$$

These quantities are rendered into individual pixels and complete RGB images at a given resolution using (neural) volume rendering.

3.8. Training Objective

We render complete views using few (n) posed input images and guide the process with k additional ground-truth views around the object. During training, we render all $n + k$ views and compare these with their ground-truth views. To evaluate the quality of the reconstruction, we apply image reconstruction losses. Concretely our reconstruction loss is:

$$\mathcal{L}(\mathbf{x}) = \frac{1}{n+k} \sum_{v=1}^{n+k} (\mathcal{L}_{\text{MSE}}(\hat{x}_v, x_v) + \alpha \mathcal{L}_{\text{LPIPS}}(\hat{x}_v, x_v)), \quad (6)$$

where \hat{x}_v is our rendered image, x_v is the ground truth image, \mathcal{L}_{MSE} is the normalized pixel-wise L2 loss, $\mathcal{L}_{\text{LPIPS}}$ is perceptual image patch similarity [43], and α is a loss weighting coefficient.

4. Experiments

In this section, we outline the experimental setup, followed by an evaluation of our implementation for two-view recon-

struction on two synthetic object datasets. We then highlight the models performance on extrapolated views, and ablate our model components.

4.1. Experimental Setup

Compared Methods. We compare the following methods: (1) Ours w/o bias, which is our model without the distance attention bias in the transformer, using only DINOv2 features and Plücker coordinate positional encodings with two input views; (2) OpenLRM [7], an open-source implementation of Large Reconstruction Models (LRM) [8] that infers a triplane representation from a single posed image, trained on a single object category; (3) PixelNeRF [41], an extension of NeRF for sparse input views that conditions on image features; and (4) Splatter Image [31], a state-of-the-art method which uses Gaussian splatting for single-view and few-view category-level reconstruction.

Datasets. We evaluate the methods on the Shapenet-SRN ‘Cars’ and ‘Chairs’ datasets [26]. Shapenet-SRN is a synthetic dataset of RGB images taken from a fixed distance around a computer model of an object, which includes camera pose information. The datasets have a predefined train/val/test split, and additional details about the datasets is included in the supplement. We use the images, camera intrinsics, camera poses and data splits as provided by the datasets and train our method using relative camera poses.

Metrics. We compare novel views with a held-out test set of ground-truth views. We measure the perceptual quality (LPIPS) [43], structural similarity (SSIM), and peak signal-to-noise ratio (PSNR) of these renders. For evaluation, we follow the protocol in other similar works using the Shapenet-SRN dataset [31, 41], by using view 64 and 128 as the input views. All unseen views are used for the computation of the metrics, excluding the conditioning input views.

Implementation and Training Details. We train an 8 layer transformer decoder model with a hidden dimension of 512. We instantiate the linear projection of our input tokens as the identity matrix, with an initial bias of zero. Each cross-attention and self-attention layer has a learnable parameter (γ) for scaling the distance bias.

Our image encoder is initialized as a pre-trained small DINOv2 vision transformer [23]. We train with $n = 2$ random input views of our instance, and $k = 2$ ground-truth comparison views (see Sec. 3.8), and compute camera poses relative to the first. This is also done at inference time. Additional training details are available in the supplement, and code is available on the project page ¹.

¹<https://github.com/SamBahrami/PluckerRF>

Table 1. Novel view synthesis results on the Chairs and Cars Shapenet-SRN test sets. Our model is given n input views, and the output renders are evaluated against all other images in the test split. Unreported results are marked with a dash. Best results are denoted in bold.

Method	n	SRN Chairs			SRN Cars		
		PSNR \uparrow	LPIPS \downarrow	SSIM \uparrow	PSNR \uparrow	LPIPS \downarrow	SSIM \uparrow
OpenLRM [7]	1	24.44	0.076	0.93	22.73	0.096	0.90
pixelNeRF [41]	1	23.72	0.128	0.90	23.17	0.146	0.89
SplatterImg [31]	1	24.43	0.067	0.93	24.00	0.078	0.92
pixelNeRF [41]	2	25.97	0.071	0.94	25.66	0.079	0.94
SplatterImg [31]	2	25.72	0.056	0.94	26.01	–	0.94
Ours w/o bias	2	27.67	0.048	0.96	25.26	0.073	0.93
Ours	2	28.22	0.045	0.96	25.54	0.070	0.94

4.2. Results

Quantitative Evaluation. We report our quantitative results for Shapenet-SRN ‘Chairs’ and ‘Cars’ in Tab. 1. Our model outperforms the baseline without the distance bias (Ours w/o bias), demonstrating that incorporating the distance bias enhances the novel view synthesis quality. Our method outperforms the 2-view implementations of Splatter Image and pixelNeRF for ‘Chairs’, and is similar in performance for the ‘Cars’ dataset.

Qualitative Results. We provide a qualitative comparison of the compared methods for the ‘Chairs’ dataset in Fig. 4 and ‘Cars’ dataset in Fig. 5, and additional qualitative results can be seen in the supplement. We can see that while our method has some blur around the unseen parts of the object, it is notably much crisper than the other methods. This is a limitation of feed-forward, diffusion-free methods in general.

Our model is challenged by areas where there is a sharp color change, such as on stripes of a car. We also note that it struggles with areas that have very fine detail such as detailed textures and geometry. These limitations may come from how the images are processed, by first converting them to lower resolution patch-wise feature maps, which may lose detail information in the process. The orthogonal feature grid is also a limiting factor since it is constrained to its grid resolution, making it challenging to capture fine details without increasing the grid size, which would have an additional memory cost.

4.3. Evaluating Extrapolated Views

In our qualitative results (Fig. 4 and Fig. 5), our novel views visually contain more high-frequency detail than the compared methods (pixelNeRF and Splatter Image), especially on the unseen side of the car. To quantify this, we re-evaluate model performance on the subset of extrapolated

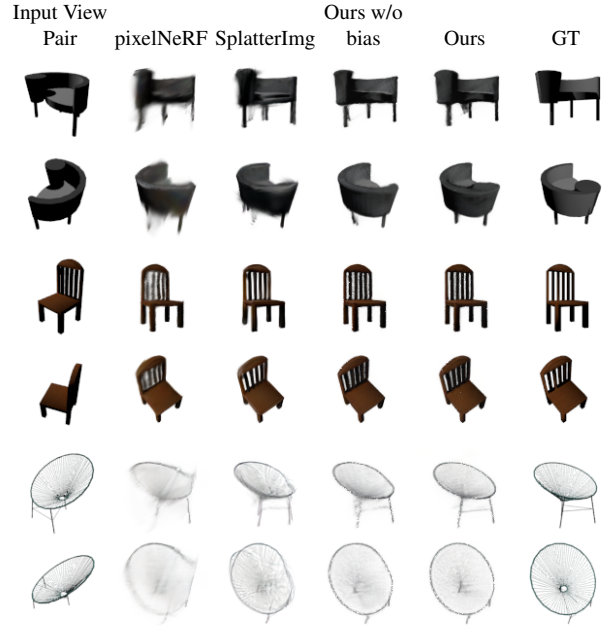


Figure 4. Shapenet-SRN Chairs qualitative comparison. The model is given two fixed input views, and renders novel views corresponding to the other ground truth views in the test set around the object. Our method is able to better infer the shape and textures of the unseen sides of the chairs, and produce sharper outputs. Please zoom in to observe finer details.

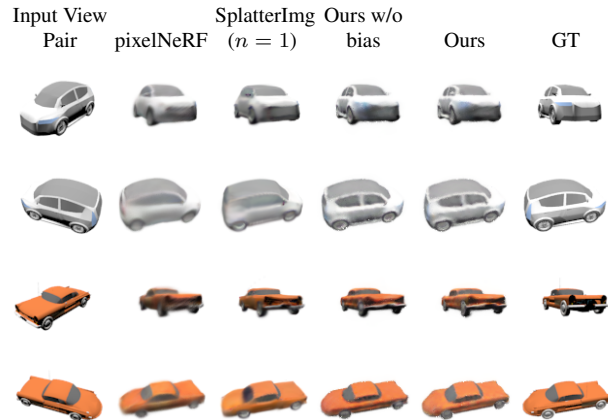


Figure 5. Shapenet-SRN Cars qualitative comparison. The model is given two fixed input views, and renders novel views corresponding to the other ground truth views in the testing set around the object. Note that the Splatter Image images in this comparison are from their single view model ($n = 1$). Please zoom in to observe finer details.

novel views, which we define as any view with an out-of-plane rotation that is at least 90° from both input views. We calculate this by computing the angle between the rotations of each input view with those of the other views in the test set, noting that the test set does not include any in-plane ro-

Table 2. Extrapolation performance. We evaluate on a subset of the Shapenet-SRN Chairs and Cars test set, corresponding to novel views that are rotated out-of-plane by at least 90° from both input views. That is, they are likely to observe the *unseen* side of the object, rather than views that interpolate the input set. We use the standard pair of input views used in the main experiments, and one view in the OpenLRM case. Bold indicates best result, unreported results are marked with a dash.

Method	SRN Chairs			SRN Cars		
	PSNR \uparrow	LPIPS \downarrow	SSIM \uparrow	PSNR \uparrow	LPIPS \downarrow	SSIM \uparrow
OpenLRM [7]	24.24	0.080	0.926	22.02	0.104	0.893
pixelNeRF [41]	25.33	0.079	0.939	23.65	0.105	0.916
SplatterImg [31]	25.15	0.064	0.936	–	–	–
Ours w/o bias	27.48	0.050	0.956	24.68	0.080	0.926
Ours	27.97	0.047	0.959	24.97	0.076	0.928

tations (i.e., rotations about the camera optical axis). The results on these extrapolated views can be seen in Tab. 2. Our method outperforms the other compared feed-forward methods more significantly under these conditions.

4.4. Ablation Study and Analysis

We conducted a series of ablation experiments to evaluate the influence of individual components of our method on the final performance. Due to computational cost, we train these models at a shorter training schedule for 60k iterations with a value of α set to 0 initially, and 0.01 for the final 20k steps on the Shapenet-SRN ‘Chairs’ dataset. Additionally we set our orthogonal feature grid resolution M to 48 instead of 64 which was used in the main experiments. All other training parameters remained consistent between these ablations and the primary experiments. We show the results of our ablation study for the two-view model in Tab. 3.

We ablate the design choice of finetuning the image encoder by instead freezing these parameters (“w/o DINOv2 finetuning”) and observe a significant decrease in the model’s performance. We ablate the choice of setting the input tokens of the transformer to a linear projection of our PlückerRF lines by instead using learnable tokens of the same size (“w/o PlückerRF input”). While this has a negligible impact on the final performance of the model, in our experiments the model converges in fewer steps with PlückerRF initialization. We ablate the choice of concatenating the Plücker coordinates to the image patch features by instead using image features only (“w/o Plücker encoding”), and observe a small drop in the model’s novel view synthesis performance. We hypothesize that the distance bias already provides most of the positional information needed by the model. We ablate the use of the perceptual similarity loss (“w/o \mathcal{L}_{LPIPS} ”) and see that the result-

Table 3. Ablation study. Experiments were undertaken on the Shapenet-SRN Chairs dataset, with two input views. The ablations are conducted on a shorter training run than the main model results, with a smaller orthogonal feature grid representation, and evaluated on the same test data split as the main experiments. Bold values indicate the highest performance in that metric, underline indicates second best.

Method	PSNR \uparrow	LPIPS \downarrow	SSIM \uparrow
Ours	24.21	0.091	0.926
w/o PlückerRF input	<u>24.17</u>	0.093	<u>0.925</u>
w/o Plücker positional encoding	23.96	0.091	0.924
w/o \mathcal{L}_{LPIPS}	24.10	0.102	0.922
w/o distance bias, Eq. (4)	23.44	0.100	0.917
w/o DINOv2 finetuning	23.10	0.101	0.913
w/o learnable γ , Eq. (4)	22.83	0.109	0.909

ing model has a similar PSNR, but also blurrier images and therefore a notable deterioration of the perceptual similarity. Finally, we ablate the learnable scaling parameters γ by setting them to 1 (“w/o learnable γ ”) and observe a significant drop in performance.

5. Discussion and Conclusion

We observe that adding our 2D–3D and 3D–3D distance biases into the attention blocks of a feed-forward few-view reconstruction transformer significantly improves its performance. However, this spatial locality bias may make it more challenging for our model to leverage symmetries in the dataset, since these would allow information to be shared non-locally. In addition, our approach is NeRF-based, which is slower to render and more memory intensive compared to other representations, such as Gaussian mixtures [12, 28, 42]. More fundamentally, feed-forward direct prediction methods are prone to blurriness in unseen regions, since they average over the possible completions, unlike diffusion-based methods that are slower but sharp.

In this paper, we introduced PlückerRF, a method for few-view 3D reconstruction that integrates 2D image information into a 3D orthogonal feature grid through a line–line distance-based attention bias. This distance bias operates (1) between input image rays and the 3D PlückerRF representation in cross-attention layers, and (2) within the 3D representation in self-attention layers, and acts to create an inductive bias in the model to help it solve the data-to-model association problem. Our model can rapidly infer a plausible 3D representation from only a few input views, outperforming existing approaches. Future work may consider a more compute-efficient 3D representation like a Gaussian mixture or introduce these inductive biases into a diffusion transformer for blur-free reconstruction.

References

- [1] Eric R. Chan, Connor Z. Lin, Matthew A. Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas Guibas, Jonathan Tremblay, Sameh Khamis, Tero Karras, and Gordon Wetzstein. Efficient geometry-aware 3D generative adversarial networks. In *CVPR*, 2022. 1, 2, 3, 6
- [2] David Charatan, Sizhe Lester Li, Andrea Tagliasacchi, and Vincent Sitzmann. pixelSplat: 3D Gaussian Splats from Image Pairs for Scalable Generalizable 3D Reconstruction. In *CVPR*, pages 19457–19467, 2024. 2
- [3] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. TensorRF: Tensorial Radiance Fields. In *ECCV*, 2022. 1, 2
- [4] Eric Ming Chen, Sidhanth Holalkere, Ruyu Yan, Kai Zhang, and Abe Davis. Ray Conditioning: Trading Photo-consistency for Photo-realism in Multi-view Image Generation. In *ICCV*, pages 23242–23251, 2023. 3, 4
- [5] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A Universe of Annotated 3D Objects. *arXiv*, 2022. 2
- [6] Jiatao Gu, Alex Trevithick, Kai-En Lin, Josh Susskind, Christian Theobalt, Lingjie Liu, and Ravi Ramamoorthi. NerfDiff: Single-image View Synthesis with NeRF-guided Distillation from 3D-aware Diffusion. In *ICML*, 2023. 2, 3
- [7] Zexin He and Tengfei Wang. OpenLrm: Open-source large reconstruction models. <https://github.com/3DTopia/OpenLRM>, 2023. 6, 7, 8, 1
- [8] Yicong Hong, Kai Zhang, Jiuxiang Gu, Sai Bi, Yang Zhou, Difan Liu, Feng Liu, Kalyan Sunkavalli, Trung Bui, and Hao Tan. LRM: Large Reconstruction Model for Single Image to 3D. In *ICLR*, 2024. 2, 3, 6, 1
- [9] Wonbong Jang and Lourdes Agapito. Codenerf: Disentangled neural radiance fields for object categories. In *CVPR*, pages 12949–12958, 2021. 2
- [10] Hanwen Jiang, Zhenyu Jiang, Kristen Grauman, and Yuke Zhu. Few-View Object Reconstruction with Unknown Categories and Camera Poses. *International Conference on 3D Vision (3DV)*, 2024. 3
- [11] Haian Jin, Hanwen Jiang, Hao Tan, Kai Zhang, Sai Bi, Tianyuan Zhang, Fujun Luan, Noah Snavely, and Zexiang Xu. LVSM: A Large View Synthesis Model with Minimal 3D Inductive Bias, 2024. 2, 3, 4
- [12] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3D Gaussian Splatting for Real-Time Radiance Field Rendering. *ACM Transactions on Graphics*, 42(4), 2023. 1, 2, 8
- [13] Alexander Kolesnikov, Alexey Dosovitskiy, Dirk Weissenborn, Georg Heigold, Jakob Uszkoreit, Lucas Beyer, Matthias Minderer, Mostafa Dehghani, Neil Houlsby, Sylvain Gelly, Thomas Unterthiner, and Xiaohua Zhai. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021. 5
- [14] Vincent Leroy, Yohann Cabon, and Jérôme Revaud. Grounding Image Matching in 3D with MAST3R. In *ECCV*, pages 71–91. Springer, 2024. 3
- [15] Jiahao Li, Hao Tan, Kai Zhang, Zexiang Xu, Fujun Luan, Yinghao Xu, Yicong Hong, Kalyan Sunkavalli, Greg Shakhnarovich, and Sai Bi. Instant3d: Fast Text-to-3D with Sparse-view Generation and Large Reconstruction Model. In *ICLR*, 2024. 2, 3
- [16] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural Sparse Voxel Fields. *NeurIPS*, 2020. 1, 2
- [17] Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick. Zero-1-to-3: Zero-shot One Image to 3D Object. In *CVPR*, pages 9298–9309, 2023. 2
- [18] Ilya Loshchilov and Frank Hutter. Decoupled Weight Decay Regularization. In *ICLR*, 2019. 1
- [19] Luke Melas-Kyriazi, Christian Rupprecht, Iro Laina, and Andrea Vedaldi. RealFusion: 360 Reconstruction of Any Object from a Single Image. In *CVPR*, 2023. 3
- [20] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis. In *ECCV*, 2020. 1, 2
- [21] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant Neural Graphics Primitives with a Multiresolution Hash Encoding. *ACM Trans. Graph.*, 41(4): 102:1–102:15, 2022. 1, 2
- [22] Norman Müller, Andrea Simonelli, Lorenzo Porzi, Samuel Rota Bulò, Matthias Nießner, and Peter Kotschieder. AutoRF: Learning 3D Object Radiance Fields from Single View Observations. In *CVPR*, 2022. 2
- [23] Maxime Oquab, Timothée Darcet, Theo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Russell Howes, Po-Yao Huang, Hu Xu, Vasu Sharma, Shang-Wen Li, Wojciech Galuba, Mike Rabbat, Mido Assran, Nicolas Ballas, Gabriel Synnaeve, Ishan Misra, Herve Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. DINOv2: Learning Robust Visual Features without Supervision. *arXiv*, 2023. 4, 6, 1
- [24] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning Continuous Signed Distance Functions for Shape Representation. In *CVPR*, 2019. 2
- [25] Julius Plücker. XVII. On a New Geometry of Space. *Philosophical Transactions of the Royal Society of London*, pages 725–791, 1865. 3
- [26] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene Representation Networks: Continuous 3D-Structure-Aware Neural Scene Representations. In *NeurIPS*, 2019. 6
- [27] Vincent Sitzmann, Semon Rezchikov, William T. Freeman, Joshua B. Tenenbaum, and Fredo Durand. Light Field Networks: Neural Scene Representations with Single-Evaluation Rendering. In *NeurIPS*, 2021. 3, 4
- [28] Brandon Smart, Chuanxia Zheng, Iro Laina, and Victor Adrian Prisacariu. Splat3r: Zero-shot Gaussian Splatting from Uncalibrated Image Pairs. 2024. 2, 3, 8

- [29] Stanislaw Szymanowicz, Christian Rupprecht, and Andrea Vedaldi. Viewset Diffusion: (0-)Image-Conditioned 3D Generative Models from 2D Data. *ICCV*, 2023. [2](#)
- [30] Stanislaw Szymanowicz, Eldar Insafutdinov, Chuanxia Zheng, Dylan Campbell, Joao Henriques, Christian Rupprecht, and Andrea Vedaldi. Flash3D: Feed-Forward Generalisable 3D Scene Reconstruction from a Single Image. *arXiv*, 2024. [2](#), [3](#)
- [31] Stanislaw Szymanowicz, Christian Rupprecht, and Andrea Vedaldi. Splatter Image: Ultra-Fast Single-View 3D Reconstruction. In *CVPR*, 2024. [2](#), [3](#), [6](#), [7](#), [8](#), [1](#)
- [32] Naveen Venkat, Mayank Agarwal, Maneesh Singh, and Shubham Tulsiani. Geometry-biased Transformers for Novel View Synthesis. *arXiv*, 2023. [2](#), [3](#)
- [33] Jianyuan Wang, Nikita Karaev, Christian Rupprecht, and David Novotny. VGGSfM: Visual geometry grounded deep structure from motion. In *CVPR*, pages 21686–21697, 2024. [3](#)
- [34] Nanyang Wang, Yinda Zhang, Zhuwen Li, Yanwei Fu, Wei Liu, and Yu-Gang Jiang. Pixel2Mesh: Generating 3D Mesh Models from Single RGB Images. In *ECCV*, 2018. [2](#)
- [35] Qianqian Wang, Yifei Zhang, Aleksander Holynski, Alexei A. Efros, and Angjoo Kanazawa. Continuous 3D Perception Model with Persistent State, 2025. [3](#)
- [36] Shuzhe Wang, Vincent Leroy, Yohann Cabon, Boris Chidlovskii, and Jerome Revaud. DUS3R: Geometric 3D Vision Made Easy. In *CVPR*, 2024. [3](#)
- [37] Olivia Wiles, Georgia Gkioxari, Richard Szeliski, and Justin Johnson. SynSin: End-to-end view synthesis from a single image. In *CVPR*, 2020. [2](#)
- [38] Rundi Wu, Yixin Zhuang, Kai Xu, Hao Zhang, and Baoquan Chen. PQ-NET: A Generative Part Seq2Seq Network for 3D Shapes. In *CVPR*, 2020. [2](#)
- [39] Jiale Xu, Weihao Cheng, Yiming Gao, Xintao Wang, Shenghua Gao, and Ying Shan. Instantmesh: Efficient 3d mesh generation from a single image with sparse-view large reconstruction models. *arXiv*, 2024. [2](#)
- [40] Zhenpei Yang, Zhile Ren, Miguel Angel Bautista, Zaiwei Zhang, Qi Shan, and Qixing Huang. FvOR: Robust Joint Shape and Pose Optimization for Few-view Object Reconstruction. In *CVPR*, pages 2497–2507, 2022. [3](#)
- [41] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelNeRF: Neural radiance fields from one or few images. In *CVPR*, 2021. [2](#), [3](#), [6](#), [7](#), [8](#), [1](#)
- [42] Kai Zhang, Sai Bi, Hao Tan, Yuanbo Xiangli, Nanxuan Zhao, Kalyan Sunkavalli, and Zexiang Xu. GS-LRM: Large Reconstruction Model for 3D Gaussian Splatting. *ECCV*, 2024. [2](#), [3](#), [8](#)
- [43] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The Unreasonable Effectiveness of Deep Features as a Perceptual Metric. In *CVPR*, 2018. [6](#)