

---

# Is Attention All NeRF Needs?

---

Mukund Varma T<sup>1†\*</sup>, Peihao Wang<sup>2†\*</sup>, Xuxi Chen<sup>2</sup>, Tianlong Chen<sup>2</sup>,  
Subhashini Venugopalan, Zhangyang Wang<sup>2†</sup>

<sup>1</sup>Indian Institute of Technology Madras, <sup>2</sup>University of Texas at Austin  
mukundvarmat@gmail.com, vsubhashini@utexas.edu  
{peihaowang, xxchen, tianlong.chen, atlaswang}@utexas.edu

## Abstract

We present *Generalizable NeRF Transformer (GNT)*, a pure, unified transformer-based architecture that efficiently reconstructs Neural Radiance Fields (NeRFs) on the fly from source views. Unlike prior works on NeRF that optimize a *per-scene* implicit representation by inverting a handcrafted rendering equation, GNT achieves *generalizable* neural scene representation and rendering, by encapsulating two transformer-based stages. The first stage of GNT, called *view transformer*, leverages multi-view geometry as an inductive bias for attention-based scene representation, and predicts coordinate-aligned features by aggregating information from epipolar lines on the neighboring views. The second stage of GNT, named *ray transformer*, renders novel views by ray marching and directly decodes the sequence of sampled point features using the attention mechanism. Our experiments demonstrate that when optimized on a single scene, GNT can successfully reconstruct NeRF without explicit rendering formula, and even improve the PSNR by  $\sim 1.3$  dB $\uparrow$  on complex scenes due to the learnable ray renderer. When trained across various scenes, GNT consistently achieves the state-of-the-art performance when transferring to forward-facing LLFF dataset (LPIPS  $\sim 20\%$   $\downarrow$ , SSIM  $\sim 25\%$   $\uparrow$ ) and synthetic blender dataset (LPIPS  $\sim 20\%$   $\downarrow$ , SSIM  $\sim 4\%$   $\uparrow$ ). In addition, we show that depth and occlusion can be inferred from the learned attention maps, which implies that *the pure attention mechanism is capable of learning a physically-grounded rendering process*. All these results bring us one step closer to the tantalizing hope of utilizing transformers as the “universal modeling tool” even for graphics. Please refer to our project page for video results: <https://vita-group.github.io/GNT/>.

## 1 Introduction

Neural Radiance Field (NeRF) [1] and its follow-up works [2, 3, 4] have achieved remarkable success on novel view synthesis, given its results on generating photo-realistic, high-resolution, and view-consistent scenes. Two key ingredients in NeRF are the coordinate-based neural network and the differentiable volumetric rendering. The coordinate-based network maps each spatial position to its corresponding color and density. Then the volume rendering pipeline will compose the points along each ray cast from the image plane to generate the target pixel color. Optimizing a NeRF can be regarded as an inverse imaging problem that overfits a neural network to satisfy the observed views. This training scheme leads to a notorious limitation of NeRF, in that it is a time-consuming optimization process for each scene [5, 6, 7].

Recent works [6, 8, 7] point out that the coordinate-based network is not necessary. They rethink novel view synthesis as a cross-view image-based interpolation problem. Unlike the vanilla NeRF that tediously fits each specific scene, those methods synthesize a generalizable 3D representation by

---

\*Equal Contribution.

†Correspondence to: Mukund Varma T, Peihao Wang and Zhangyang Wang.

aggregating image features extracted from seen views according to the camera and geometry priors [7, 6, 1]. Despite notable improvements achieved by these line of works, they all decode the feature volume to a radiance field, and rely on the classical volume rendering [9] to generate images. Note that, volume rendering was initially proposed to handle translucent and non-surface materials [10]. It lacks optical modeling of a solid surface [11, 12], reflectance [13, 14, 4], inter-surface scattering and else. This implies that radiance fields along with volume rendering are not a universal imaging model, which might have limited the generalization ability of NeRFs as well.

In this paper, we first consider the problem of transferable novel view synthesis as a two-stage information aggregation process: the multi-view image feature fusion, followed by the sampling-based rendering integration. It has not escaped from our notice that transformer [15], with ubiquitous success in computer vision [16, 17, 18, 19], processing sequential data by passing information globally across tokens, can be served a universal aggregation function [20]. To this end, we propose to model these two stages with purely attention-based transformer architectures. For volumetric scene representation, we adopt the first transformer, dubbed *view transformer*, to aggregate pixel-aligned image features [1] from corresponding epipolar lines to predict coordinate-wise features. When rendering a novel view, we employ the second transformer, dubbed *ray transformer* that performs ray tracing and composes point features along the ray via attention mechanism. These two parts together assemble our solution dubbed *Generalizable NeRF Transformer (GNT)*.

Remarkably, GNT does not require any explicit modeling of the scene representation nor the rendering pipeline. Instead, it simultaneously learns to represent scenes from source view images and perform a learnable and scene adaptive ray-based rendering using the attention mechanism. Therefore, GNT can predict novel views using the captured images without per-scene fitting. Our promising results endorses that transformers are strong, scalable and a versatile learning backbone even for graphical rendering [21]. Our technical contributions can be summarized as follows:

- We propose a purely transformer-based NeRF architecture, dubbed GNT, that achieves more expressive and universal scene representation and rendering by unifying coordinate networks and volumetric renderer into a two-stage transformer, being capable of generalizing to unseen scenes when trained across instances.
- For a more expressive volumetric representation, GNT employs the *view transformer* to aggregate multi-view image features complying with epipolar geometry to infer coordinate-aligned features. To learn a more universal ray-based rendering, GNT utilizes the *ray transformer* to predict the ray color. Together these two parts assemble a transformer that renders novel views by completely relying on the attention mechanism, and inherently learns to be depth- and occlusion-aware.
- We empirically demonstrate that GNT significantly improves the PSNR of existing NeRFs (single scene) by up to  $\sim 1.3$  dB in complex scenes. In the cross-scene generalization scenario, GNT achieves state-of-the-art perceptual metric scores by outperforming other baselines by up to  $\sim 20\%$   $\downarrow$  LPIPS and  $\sim 12\%$   $\uparrow$  SSIM.

## 2 Related Work

**Advances in Transformer.** Transformer [15] have emerged as a novel learning backbone that utilizes self-attention mechanism to effectively capture long-range correlation for sequential data. Undebatable success has been witnessed in the domain of natural language processing [22, 23, 24], computer vision [16, 25, 18], and cross-disciplinary subjects [26, 27, 28]. In computer vision, [16] established the first Vision Transformer (ViT) for image classification task and achieved state-of-the-art results on ImageNet. In architectural design, [17, 29] further improve its data efficiency and scalability. [25] proposes the window-based attention to leverage multi-scale information and reduce computational complexity. In application level, other follow-up works extended ViT to various classic vision tasks, such as object detection [30, 31, 32, 33], segmentation [34, 35], video processing [36, 37], and 3D instance processing [38, 39]. [18] also shows that transformer is even capable of image generation. However, conventional image generation mainly rely on a statistical model lacking 3D priors, and thus suffers from view inconsistency and low-resolution.

**Neural Radiance Fields.** Neural Radiance Field (NeRF) [40] achieves photorealistic and consistent novel view synthesis results by fitting each scene as a continuous 5D radiance field parameterized by

an MLP. Mip-NeRF [2, 41] addresses sampling and aliasing problems for NeRF. Recent works of [4, 12, 11, 42] improve the surface representation, [43, 44, 45] extend NeRF to dynamic scenes, and [13, 14] introduce lighting and reflection modeling to NeRF. To avoid per-scene training, PixelNeRF [7], IBRNet [6], MVNeRF [5], and NeuRay [46] reconstruct the radiance field with one-shot forward pass via training a cross-scene multi-view aggregator. [47, 48, 49, 50, 51] accelerate NeRF in both training and inference.

**Transformer Meets Radiance Fields.** Several recent attempts have been made to apply transformers for novel view synthesis. For cross-scene generalization, IBRNet [6] adopts a transformer to pass information among sampled points on the ray to predict density. NerFormer [52] uses attention module to aggregate source views to construct feature volume with epipolar geometry constraints. [53] proposes similar transformer-based architecture to encode multi-view information. These three works demonstrate the promise of transformer as a core component in NeRF. However, all of them decode the latent feature representation to point-wise color and density, and rely on classic volume rendering to form image, which ignores transformer’s potential in rendering. NLF [8] establishes an attention based framework to display light fields with view consistency, where the first transformer summarizes information on epipolar lines independently and then fuse epipolar features using the second transformer. We point out that their aggregation strategy differs from our GNT, and fails to generalize across scenes due to the lacking of cross-view communication. [54] divides a single image into patches and converts them into 3D representation via a vision transformer. SRT [55] treats images and camera parameters all in some latent space, and trains a transformer that directly maps camera pose embedding to the corresponding image without any physical constraints. This approach presumably requires heavy training resources and is difficult to scale up to higher resolution.

### 3 Preliminaries

**Self-Attention and Transformer.** Transformer-based architecture tokenizes data into sequences and purely relies on attention mechanism to process them, which demonstrates promise in both language [15] and vision [16] domains. Multi-Head Self-Attention (MHA), the main ingredient of transformer, exchanges information amongst tokens according to a pairwise computed score, which indicates the focus on a subset of context. Formally, let  $\mathbf{X} \in \mathbb{R}^{N \times d}$  represent some sequential data with  $N$  tokens of  $d$ -dimension. A self-attention layer transforms the feature matrix as below:

$$\text{Attn}(\mathbf{X}) = \text{softmax}(\mathbf{A})f_V(\mathbf{X}), \text{ where } \mathbf{A}_{i,j} = \alpha(\mathbf{X}_i, \mathbf{X}_j), \forall i, j \in [N] \quad (1)$$

where  $\mathbf{A} \in \mathbb{R}^{N \times N}$  is called the attention matrix,  $\text{softmax}(\cdot)$  operation normalizes the attention matrix row-wisely, and  $\alpha(\cdot)$  represents a pair-wise relation function. The most commonly adopted relation mapping is dot-product  $\alpha(\mathbf{X}_i, \mathbf{X}_j) = f_Q(\mathbf{X}_i)^\top f_K(\mathbf{X}_j)/\gamma$ , where  $f_Q(\cdot), f_K(\cdot), f_V(\cdot)$  are called query, key, and value mapping functions. In a standard transformer, they are chosen as fully-connected layers. The functionality of self-attention can be regarded as an aggregation operator, which pass information from one token to another token according to the attention value between them. Multi-Head Self-Attention (MHA) sets a group of self-attention blocks, and adopts a linear layer to project them onto the output space:

$$\text{MHA}(\mathbf{X}) = [\text{Attn}_1(\mathbf{X}) \quad \text{Attn}_2(\mathbf{X}) \quad \cdots \quad \text{Attn}_H(\mathbf{X})] \mathbf{W}_O \quad (2)$$

Following an MHA block, one standard layer of transformer also adopts a Feed-Forward Network (FFN) to do point-wise feature transformation, as well as skip connection and layer normalization to stabilize training. The whole transformer block can be formulated as below:

$$\tilde{\mathbf{X}} = \text{MHA}(\text{LayerNorm}(\mathbf{X})) + \mathbf{X}, \quad \mathbf{Y} = \text{FFN}(\text{LayerNorm}(\tilde{\mathbf{X}})) + \tilde{\mathbf{X}} \quad (3)$$

**Neural Radiance Field.** Neural Radiance Fields (NeRFs) [40] converts multi-view images into a radiance field and interpolates novel views by re-rendering the radiance field from a new angle. Technically, NeRF models the underlying 3D scene as a continuous radiance field  $\mathcal{F} : (\mathbf{x}, \boldsymbol{\theta}) \mapsto (\mathbf{c}, \sigma)$  parameterized by an Multi-Layer Perceptron (MLP)  $\Theta$ , which maps a spatial coordinate  $\mathbf{x} \in \mathbb{R}^3$  together with the viewing direction  $\boldsymbol{\theta} \in [-\pi, \pi]^2$  to a color  $\mathbf{c} \in \mathbb{R}^3$  plus density  $\sigma \in \mathbb{R}_+$  tuple. To form an image, NeRF performs the ray-based rendering, where it casts a ray  $\mathbf{r} = (\mathbf{o}, \mathbf{d})$  from the optical center  $\mathbf{o} \in \mathbb{R}^3$  through each pixel (towards direction  $\mathbf{d} \in \mathbb{R}^3$ ), and then leverages volume

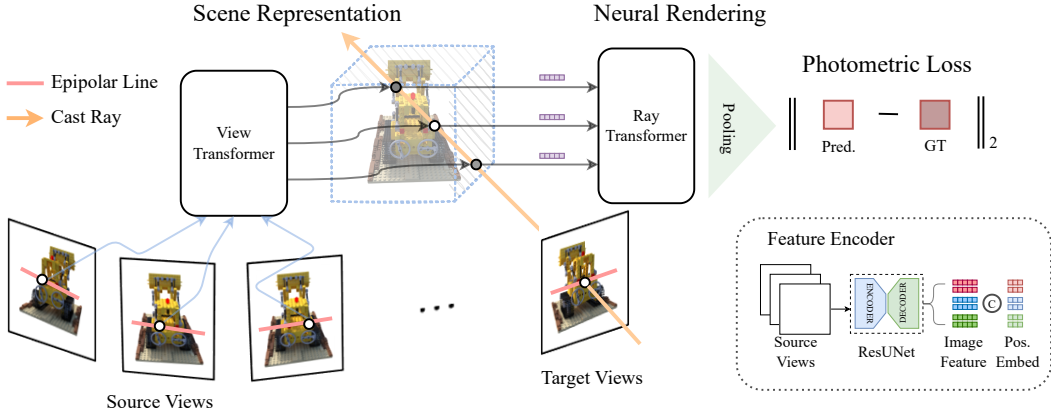


Figure 1: Overview of GNT: 1) Identify source views for a given target view, 2) Extract features for epipolar points using a trainable U-Net like model, 3) For each ray in the target view, sample points and directly predict target pixel’s color by aggregating view-wise features (View Transformer) and across points along a ray (Ray Transformer).

rendering [9] to compose the color and density along the ray between the near-far planes:

$$C(\mathbf{r}|\Theta) = \int_{t_n}^{t_f} T(t)\sigma(\mathbf{r}(t))\mathbf{c}(\mathbf{r}(t), \mathbf{d})dt, \text{ where } T(t) = \exp\left(-\int_{t_n}^t \sigma(s)ds\right), \quad (4)$$

where  $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$ ,  $t_n$  and  $t_f$  are the near and far planes, respectively. In practice, the Eq. 4 is numerically estimated using quadrature rules [40]. Given images captured from surrounding views with known camera parameters, NeRF fits the radiance field by maximizing the likelihood of simulated results. Suppose we collect all pairs of rays and pixel colors as the training set  $\mathcal{D} = \{(\mathbf{r}_i, \widehat{\mathbf{C}}_i)\}_{i=1}^N$ , where  $N$  is the total number of rays sampled, and  $\widehat{\mathbf{C}}_i$  denotes the ground-truth color of the  $i$ -th ray, then we train the implicit representation  $\Theta$  via the following loss function:

$$\mathcal{L}(\Theta|\mathcal{R}) = \mathbb{E}_{(\mathbf{r}, \widehat{\mathbf{C}}) \in \mathcal{D}} \|C(\mathbf{r}|\Theta) - \widehat{\mathbf{C}}(\mathbf{r})\|_2^2, \quad (5)$$

## 4 Method: Make Attention All NeRF Needs

**Overview.** Given a set of  $N$  input views with known camera parameters  $\{(\mathbf{I}_i \in \mathbb{R}^{H \times W \times 3}, \mathbf{P}_i \in \mathbb{R}^{3 \times 4})\}_{i=1}^N$ , our goal is to synthesize novel views from arbitrary angles. In contrast to vanilla NeRF which is limited by per-scene fitting, we follow the design of cross-scene generalizable NeRFs [7, 5, 6]. Our method can be divided into two stages: 1) construct the 3D representation from source views on the fly in the feature space, 2) re-render the feature field at the specified angle to synthesize novel views. However, unlike [7, 5, 6, 46], which design dedicated multi-view aggregators for scene representation and borrow classic volume rendering for view synthesis, we propose a unified transformer-based architecture, dubbed GNT, to model both stages. Our pipeline is depicted in Fig. 1. For the first stage, we propose the *view transformer* to aggregate coordinate-aligned features from source views. To enforce multi-view geometry, we inject the inductive bias of epipolar constraints into the attention mechanism. After obtaining the feature representation of each point on the ray, GNT leverages the *ray transformer* to compose point-wise features along the ray to form the ray color. The whole pipeline is purely attention based and can be trained end-to-end.

### 4.1 Epipolar Geometry Constrained Scene Representation

**Multi-View Feature Aggregation.** NeRF represents 3D scene as a radiance field  $\mathcal{F} : (\mathbf{x}, \theta) \mapsto (\mathbf{c}, \sigma)$ , where each point corresponds to a color and density tuple. Vanilla NeRF [40] parameterizes the radiance field using an MLP, and recovers the scene in a backward optimization fashion, which inherently hinders NeRF from generalizing to new scenes. In contrast, generalizable NeRFs [7, 6, 5] construct the radiance field in a feed-forward scheme, which directly encodes multi-view images into

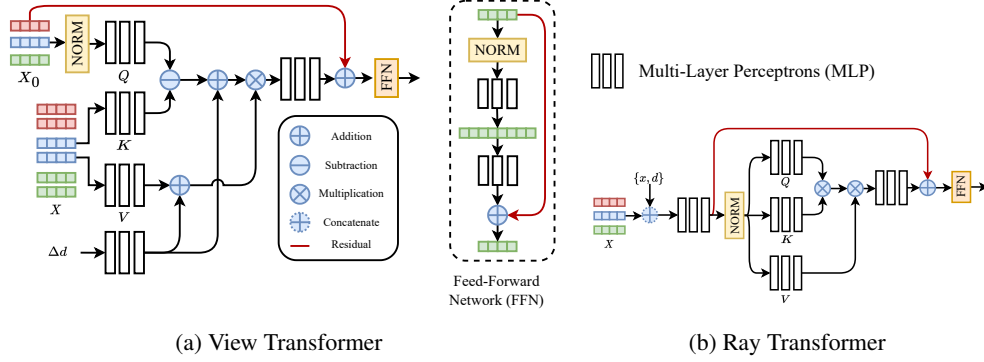


Figure 2: Detailed network architectures of view transformer and ray transformer in GNT.

3D feature space, and decodes it to a color-density field. Instead of transforming 3D features into physical variables (e.g., color and density), we implicitly represent a 3D scene as a coordinate-aligned feature field  $\mathcal{F} : (\mathbf{x}, \boldsymbol{\theta}) \mapsto \mathbf{f} \in \mathbb{R}^d$ , where  $d$  is the dimension of the latent space. We formulate the feed-forward scene representation as follows:

$$\mathcal{F}(\mathbf{x}, \boldsymbol{\theta}) = \mathcal{V}(\mathbf{x}, \boldsymbol{\theta}; \{\mathbf{I}_1, \dots, \mathbf{I}_N\}), \quad (6)$$

where  $\mathcal{V}(\cdot)$  is a function invariant to the permutation of input images to aggregate different views  $\{\mathbf{I}_1, \dots, \mathbf{I}_N\}$  into a coordinate-aligned feature field, and extracts features at a specific location. It has been shown that transformers can serve as a universal set aggregation function [20]. However, naively plugging in the attention mechanism that globally attends to every pixel in the source images [55, 56] is memory prohibitive and lacks multi-view geometric priors. To this end, we introduce epipolar geometry as an inductive bias that restricts each pixel to only attend to pixels that lie on the corresponding epipolar lines of the neighboring views. We name this epipolar geometry based cross-view aggregation transformer  $\mathcal{V}(\cdot)$  as *view transformer*. Specifically, we divide view transformer into the following steps. First we extract a feature map  $\mathbf{F}_i = \text{UNet}(\mathbf{I}_i) \in \mathbb{R}^{H \times W \times d}$  from each image via a shared U-Net based convolutional neural network [57, 6]. To obtain the feature representation at a position  $\mathbf{x}$ , we first project  $\mathbf{x}$  to every source image, and interpolate the feature vector on the image plane. Afterwards, we adopt a transformer to combine all the feature vectors. Formally, this process can be written as below:

$$\mathcal{V}(\mathbf{x}, \boldsymbol{\theta}; \{\mathbf{I}_1, \dots, \mathbf{I}_N\}) = \text{View-Transformer}(\mathbf{F}_1(\Pi_1(\mathbf{x}), \boldsymbol{\theta}), \dots, \mathbf{F}_N(\Pi_N(\mathbf{x}), \boldsymbol{\theta})), \quad (7)$$

where  $\text{View-Transformer}(\cdot)$  is a multi-layer stacking of Eq. 3.  $\Pi_i(\mathbf{x}) = [x/z \ y/z]^\top$  with  $[x \ y \ z]^\top = \mathbf{P}_i [\mathbf{x} \ 1]^\top$  projects 3D point  $\mathbf{x} \in \mathbb{R}^3$  onto the  $i$ -th image plane.  $\mathbf{F}_i(\mathbf{x}', \boldsymbol{\theta}) \in \mathbb{R}^d$  computes the feature vector at position  $\mathbf{x}' \in \mathbb{R}^2$  via bilinear interpolation on the feature grids. We also concatenate the extracted feature vector with point coordinate  $\gamma(\mathbf{x})$ , viewing direction  $\gamma(\boldsymbol{\theta})$ , and relative directions of source views with respect to the target view  $\Delta \mathbf{d}_i = \gamma(\mathbf{x} - \mathbf{o}_i)$  [6], where  $\mathbf{o}_i$  is the origin point of the  $i$ -th camera, and  $\gamma(\cdot)$  denotes a position encoding function.

**Memory-Efficient Cross-View Attention.** Computing attention between every pair of inputs has  $O(N^2)$  memory complexity, which is computational prohibitive when sampling thousands of points at the same time. Nevertheless, we note that view transformer only needs to read out one token as the fused results of all the views. Therefore, we propose to only place one read-out token  $\mathbf{X}_0 \in \mathbb{R}^d$  in the query sequence, and let it iteratively summarize features from other data points. This reduces the complexity for each layer up to  $O(N)$ . We initialize the read-out token as the element-wise max-pooling of all the inputs:  $\mathbf{X}_0 = \max(\mathbf{F}_1(\Pi_1(\mathbf{x}), \boldsymbol{\theta}), \dots, \mathbf{F}_N(\Pi_N(\mathbf{x}), \boldsymbol{\theta}))$ . Rather than adopting a standard dot-product attention, we choose subtraction operation as the relation function. Subtraction attention has been shown more effective for positional and geometric relationship reasoning [58, 59]. Compared with dot-product that collapses the feature dimension into a scalar, subtraction attention computes different attention scores for every channel of the value matrix, which increases diversity in feature interactions. Moreover, we augment the attention map and value matrix with  $\{\Delta \mathbf{d}_i\}_{i=1}^N$  to provide relative spatial context. Technically, we utilize a linear layer  $\mathbf{W}_P$  to lift  $\Delta \mathbf{d}_i$  to the hidden dimension. We illustrate view transformer in Fig. 2a. To be more specific, the modified attention

adopted in our view transformer can be formulated as:

$$\text{View-Attn}(\mathbf{X}) = \text{diag}(\text{softmax}(\mathbf{A} + \mathbf{\Delta}^\top) f_V(\mathbf{X} + \mathbf{\Delta})), \text{ where } \mathbf{A}_j = f_Q(\mathbf{X}_0) - f_K(\mathbf{X}_j), \quad (8)$$

where  $\mathbf{A}_j \in \mathbb{R}^d$  denotes the  $j$ -th column of  $\mathbf{A}$ ,  $\mathbf{\Delta} = [\mathbf{\Delta}d_1 \ \cdots \ \mathbf{\Delta}d_N]^\top \mathbf{W}_P \in \mathbb{R}^{N \times d}$ ,  $f_Q$ ,  $f_K$ , and  $f_V$  are parameterized by an MLP. We note that by applying  $\text{diag}(\cdot)$ , we read out the updated query token  $\mathbf{X}_0$ . An efficient implementation is presented in Appendix A.

**Discussion on Occlusion Awareness.** Conceptually, view transformer attempts to find correspondence between the queried points and source views. The learned attention amounts to a likelihood score that a pixel on the source view is an image of the same point in the 3D space, i.e., no points lies between the target point and the pixel. NeuRay [46] leverages the cost volume from MVSNet [60] to predict per-pixel visibility and shows that introducing occlusion information is beneficial for multi-view aggregation in generalizable NeRF. We argue that instead of explicitly regressing the visibility, purely relying on epipolar geometry constrained attention can automatically learn how to infer occlusion, given prior works in Multi-View Stereo (MVS) [61, 19]. In view transformer, the U-Net provides multi-scale features to the transformer, and the attention block acts as a matching algorithm, which selects the pixels from neighboring views that maximize view consistency. We defer more discussion to Sec. 5.5.

## 4.2 Attention Driven Volumetric Rendering

Volume rendering (Eq. 4) has been regarded as a key knob of NeRF’s success, which simulates light transport and occlusion in a radiance field. Generalizable NeRFs [7, 6, 5, 46] also inherit the rendering function from NeRF and explicitly decode intermediate 3D latent representation to color-density field. However, volume rendering still struggles to model sharp surfaces [42, 11, 12] and complex interference patterns, such as specular reflection [13, 14], refraction, and translucency [8]. Even though physical imaging modeling have been investigated in prior literature, none of them can unify a rendering equation that universally applies to all scenarios. This motivates us to replace the handcrafted rendering function with a data-driven renderer that learns to project a 3D feature field onto 2D images with respect to specified camera poses.

**Ray Feature Aggregation.** To illustrate how we generalize NeRF’s volume rendering, we first follow the ray-based rendering pipeline proposed in [40]. That being said, we first cast a ray  $\mathbf{r} = (\mathbf{o}, \mathbf{d})$  from the camera origin through the pixel to be rendered, and then compute the numerical estimation of Eq. 4 by randomly sampling  $M$  points  $\{\mathbf{x}_i = \mathbf{o} + t_i \mathbf{d}\}_{i=1}^M$  along the ray. We rewrite the numerical estimation of Eq. 4 as a weighted summation:

$$\mathbf{C}(\mathbf{r}) = \sum_{i=1}^M T(i)(1 - \exp(-\sigma_i \delta_i)) \mathbf{c}_i := \sum_{i=1}^M w_i \mathbf{c}_i, \text{ where } T(i) = \exp\left(-\sum_{j=1}^{i-1} \sigma_j \delta_j\right), \quad (9)$$

and  $w_i = T(i)(1 - \exp(-\sigma_i \delta_i))$ . One can see volume rendering is a weighted aggregation of all the point-wise output, in which the weights are dependent on the other points for occlusion modeling. We notice that this process shares the intrinsic similarity with the computational paradigm of the attention, where token features correspond to  $\mathbf{c}_i$ , and attention scores correspond to  $w_i$ . Therefore, we can generalize Eq. 9 with a *ray transformer*, whose one layer is illustrated in Fig. 2b. Following view transformer, we can compute a feature representation  $\mathbf{f}_i = \mathcal{F}(\mathbf{x}_i, \boldsymbol{\theta}) \in \mathbb{R}^d$  for each sampled point  $\mathbf{x}_i$ . In addition to this, we also add spatial location  $\gamma(\mathbf{x})$  and view direction  $\gamma(\boldsymbol{\theta})$  into  $\mathbf{f}_i$ . We obtain the rendered ray color by feeding the sequence of  $\{\mathbf{f}_1, \cdots, \mathbf{f}_M\}$  into the ray transformer, perform mean pooling over all the predicted tokens, and map the pooled feature vector to RGB via an MLP. The whole process can be formulated as:

$$\mathbf{C}(\mathbf{r}) = \text{MLP} \circ \text{Mean} \circ \text{Ray-Transformer}(\mathcal{F}(\mathbf{o} + t_1 \mathbf{d}, \boldsymbol{\theta}), \cdots, \mathcal{F}(\mathbf{o} + t_M \mathbf{d}, \boldsymbol{\theta})), \quad (10)$$

where  $t_i \sim \mathcal{U}[t_n + (i - 1)(t_f - t_n)/M, t_n + i(t_f - t_n)/M]$ . Unlike view transformer, we adopt the standard dot-product attention in ray transformer, which encourages more cross-point interaction to provide enough contextual information for accurate occupancy and occlusion prediction (which resembles plane-sweep stereo methods considering matching scores over all hypothesis planes [6]). We argue that ray transformer can automatically learn the attention distribution to control the

Table 1: Comparison of GNT against SOTA methods for single scene rendering on the LLFF dataset.

Models	Orchids				Trex				Horns			
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	Avg $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	Avg $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	Avg $\downarrow$
LLFF [62]	18.52	0.588	0.141	0.313	24.15	0.857	0.222	0.068	24.70	0.840	0.193	0.064
NeRF [40]	20.36	0.641	0.321	0.121	26.80	0.880	0.249	0.056	27.45	0.828	0.268	0.058
NeX [63]	20.42	0.765	0.242	0.102	28.73	0.953	0.192	0.038	28.46	0.934	0.173	0.040
NLF [8]	21.05	<b>0.807</b>	0.173	0.084	<b>30.34</b>	<b>0.968</b>	0.143	<b>0.029</b>	<b>29.78</b>	<b>0.957</b>	0.121	0.030
GNT	<b>21.20</b>	0.760	<b>0.161</b>	<b>0.084</b>	27.62	0.928	<b>0.091</b>	0.034	29.41	0.931	<b>0.085</b>	<b>0.029</b>

sharpness of reconstructed surface, and bake desirable lighting effects from the illumination and material aware features extracted from source images. The pseudocode of our implementation is provided in Appendix A. We also involve discussion on the extension to auto-regressive rendering and attention-based coarse-to-fine sampling in Appendix B.

**Discussion on Depth Cuing.** The ray transformer iteratively aggregates features according to the attention value. This attention value can be regarded as the importance of each point to form the image, which reflects visibility and occlusion reasoned by point-to-point interaction. Therefore, we can interpret the average attention score for each point as the accumulated weights  $w_i$  in Eq. 9. In this sense, we can infer the depth map from the attention map by averaging the marching distance  $t_i$  with the attention value. This implies our ray transformer learns geometry-aware 3D semantics on both feature space and attention map, which helps it generalize well across scenes. We defer visualization and analysis to Sec. 5.5. NLF [8] proposes a similar two-stage rendering transformer, but it first extracts features on the epipolar lines and then aggregates epipolar features to get the pixel color. We doubt this strategy may fail to generalize as epipolar features lack communication with each other and thus cannot induce geometry grounded semantics.

## 5 Experiments

We conduct extensive experiments to compare GNT with several state-of-the-art (SOTA) methods for novel view synthesis. We first provide qualitative and quantitative results in the single scene training setting, followed by extending our approach for generalization to unseen scenes.

### 5.1 Implementation Details

**Source and Target view sampling.** As described in [6], we construct a training pair of source and target views by first selecting a target view, then identifying a pool of  $k \times N$  nearby views, from which  $N$  views are randomly sampled as source views. This sampling strategy simulates various view densities during training and therefore helps the network generalize better. During training, the values for  $k$  and  $N$  are uniformly sampled at random from [1, 3] and [8, 12] respectively.

**Network Architecture.** To extract features from the source views, we use a U-Net like architecture with a ResNet34 encoder, followed by two up-sampling layers as decoder [57]. Each view transformer block contains a single-headed cross-attention layer while the ray transformer block contains a multi-headed self-attention layer with four heads. The outputs from these attention layers are passed onto corresponding feedforward blocks with a Rectified Linear Unit (RELU) activation and a hidden dimension of 256. A residual connection is applied between the pre-normalized inputs (LayerNorm) and outputs at each layer. For all our single scene experiments, we alternatively stack 4 view and ray transformer blocks while our larger generalization experiments use 8 blocks each. All transformer blocks (view and ray) are of dimension 64. Following [15, 40, 64], we convert the low-dimensional coordinates to a high-dimensional representation using Fourier components, where the number of frequencies is selected as 10 for all our experiments. The derived view and position embeddings are each of dimension 63.

**Training / Inference Details.** We train both the feature extraction network and GNT end-to-end on datasets of multi-view posed images using the Adam optimizer to minimize the mean-squared error between predicted and ground truth RGB pixel values. The base learning rates for the feature extraction network and GNT are  $10^{-3}$  and  $5 \times 10^{-4}$  respectively, which decay exponentially over

Table 2: Comparison of GNT against SOTA methods for single scene rendering on the Blender dataset.

Models	Lego				Chair				Drums			
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	Avg $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	Avg $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	Avg $\downarrow$
LLFF [62]	24.54	0.911	0.110	0.049	28.72	0.948	0.064	0.027	21.13	0.890	0.126	0.069
NeRF [40]	32.54	0.961	0.050	0.018	33.00	0.967	0.046	0.016	25.01	0.925	0.091	0.043
MipNeRF [2]	35.7	0.978	0.021	0.009	35.14	0.981	0.021	0.010	25.48	0.932	0.065	0.036
NLF [8]	<b>35.76</b>	<b>0.989</b>	<b>0.010</b>	<b>0.007</b>	<b>35.30</b>	<b>0.989</b>	<b>0.012</b>	<b>0.007</b>	25.83	0.955	<b>0.045</b>	0.029
GNT	31.43	0.967	0.032	0.016	33.10	0.979	0.023	0.010	<b>28.43</b>	<b>0.957</b>	0.046	<b>0.023</b>

training steps. For all our experiments, we train for 250,000 steps with 2048 rays sampled in each iteration. Unlike most NeRF methods, we do not use separate coarse, fine networks and sample only 64 coarse points per ray across all experiments (unless otherwise specified).

**Metrics.** To measure the performance of our model we use three widely adopted metrics: Peak Signal-to-Noise Ratio (PSNR), Structural Similarity Index Measure (SSIM) [65], and the Learned Perceptual Image Patch Similarity (LPIPS) [66]. We report the averages of each metric over different views in each scene (for single scene experiments) and across multiple scenes in each dataset (for generalization experiments). We additionally report the geometric mean of  $10^{-PSNR/10}$ ,  $\sqrt{1 - SSIM}$  and LPIPS, which provides a summary of the three metrics for easier comparison [2].

## 5.2 Single Scene Results.

**Datasets.** To evaluate the single scene view generation capacity of GNT, we perform experiments on datasets containing synthetic rendering of objects and real images of complex scenes. In these experiments, we use the same resolution and train/test splits as NeRF [40].

(a) **Local Light Field Fusion (LLFF) dataset:** The LLFF dataset introduced by [62] consists of 8 forward facing captures of real-world scenes using a smartphone. Due to resource constraints, we select three scenes - Trex, Horns, Orchids and the computed metrics are summarized in Tab. 1.

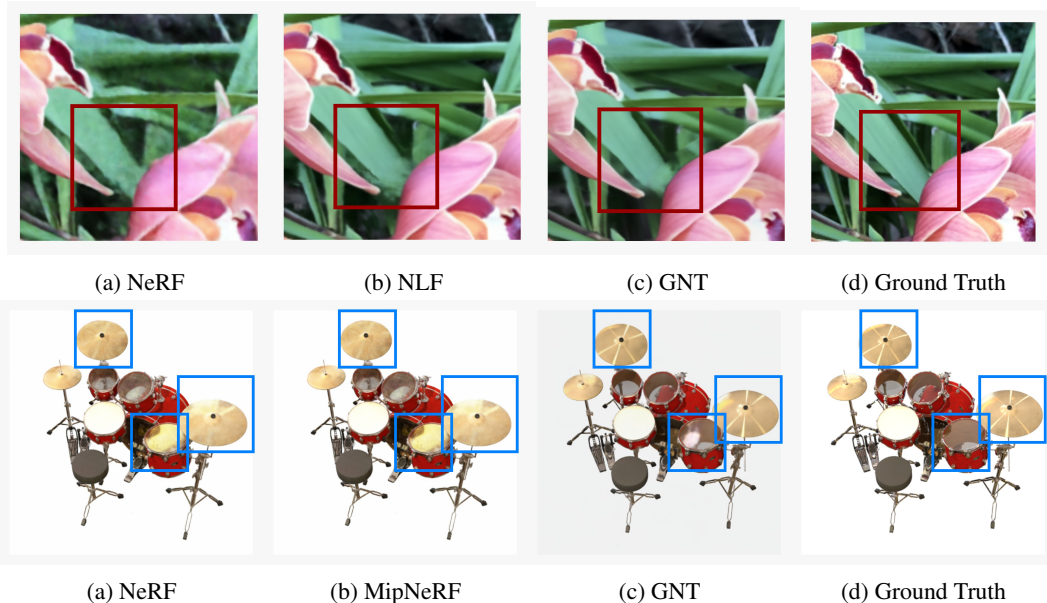


Figure 3: Qualitative results of GNT for single-scene rendering. In the Orchids scene from LLFF dataset (first row), GNT recovers the shape of the leaves more accurately. In the Drums scene from Blender dataset (second row), GNT’s learnable renderer is able to model physical phenomenon like specular reflections.



Table 3: Comparison of GNT against SOTA methods for cross scene generalization on the LLFF, Blender datasets

Models	Real Forward-Facing				Blender			
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	Avg $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	Avg $\downarrow$
PixelNeRF [7]	18.66	0.588	0.463	0.159	22.65	0.808	0.202	0.078
MVSNeRF [5]	21.18	0.691	0.301	0.108	25.15	0.853	0.159	0.057
IBRNet [6]	25.17	0.813	0.200	0.064	26.73	0.908	0.101	0.040
NeuRay [46]	<b>25.35</b>	0.818	0.198	0.062	<b>28.29</b>	0.927	0.080	0.032
GNT	25.18	<b>0.838</b>	<b>0.153</b>	<b>0.057</b>	26.94	<b>0.931</b>	<b>0.062</b>	<b>0.032</b>

(b) **Blender dataset:** The synthetic dataset introduced by [40] consists of 8, 360° scenes of objects with complicated geometry and realistic material. Each scene consists of images rendered from viewpoints randomly sampled on a hemisphere around the object. Similar to the experiments on the LLFF dataset, we select three scenes - Lego, Chair, Drums, and report metrics in Tab. 2.

**Discussion.** We compare our GNT with LLFF [62], NeRF [40], MipNeRF [2], NeX [63], and NLP [8]. Compared to other methods, we utilize a smaller batch size and sample fewer points per ray during training. Specifically, GNT samples 2048 rays per batch and 64 coarse points per ray while NLF samples as much as 8192 rays per batch and 64 coarse, 128 fine points per ray. These hyperparameters have a strong correlation with the rendered image quality, leaving our method at a disadvantage. However, a network trained only on 64 points can be evaluated using more points due to the input length independent property of the attention operation. Therefore, to bring GNT to a comparable experimental setup with other methods, we present results on images rendered by GNT with increased sampling points (i.e., 192 points per ray). Despite these differences, GNT still manages to outperform existing SOTA methods on the LLFF dataset by  $\sim 13\%$   $\downarrow$  average score, and on-par in the Blender dataset. In complex scenes, GNT manages to outperform other methods more substantially by 0.15 dB, 1.6 dB in the Orchids and Drums scene respectively. This indicates the effectiveness of our attention-driven volumetric rendering to model complex conditions. Interestingly, even in the worse performing scenes by PSNR (e.g. Trex), GNT achieves best perceptual metric scores (LPIPS  $\sim 20\%$   $\downarrow$ , SSIM  $\sim 25\%$   $\uparrow$ ). This could be because PSNR fails to measure structural distortions, blurring, has high sensitivity towards brightness, and hence does not effectively measure visual quality. Similar inferences are discussed in [67] regarding discrepancies in PSNR scores and their correlation to rendered image quality. In the Blender dataset, we do not observe similar performance improvements across all scenes as the views are not dense enough to cover the entire 360°, therefore not complying with the epipolar constraints. As discussed in Sec. 4.2, the attention scores obtained from the ray transformer can be interpreted as density values of each point. Due to the softmax normalization step, the sum of density values (or occupancy) would always be equal to 1, which is not valid in scenes with empty backgrounds, e.g. Blender. Therefore, we find the error rate to be maximum in regions corresponding to such empty backgrounds and could further possibly explain the relatively weaker performance of GNT in the Blender dataset. Figs. 3 provides qualitative comparisons on the Orchids, Drums scene respectively and we can clearly see that GNT recovers the edge details of objects (in the case of Orchids) and models complex lighting effect like specular reflection (in the case of Drums) more accurately.

### 5.3 Generalization to Unseen Scenes.

**Datasets.** GNT leverages multi-view features complying with epipolar geometry, enabling generalization to unseen scenes. We follow the experimental protocol described in [6] to evaluate the cross-scene generalization capacity of GNT and use the following datasets for training, evaluation respectively.

(a) **Training Datasets:** The training dataset consists of both real, and synthetic data. For synthetic data, we use object renderings of 1023 models from Google Scanned Object [68]. For real data, we make use of RealEstate10K [69], 100 scenes from the Spaces dataset [70], and 102 real scenes from handheld cellphone captures [62, 6].

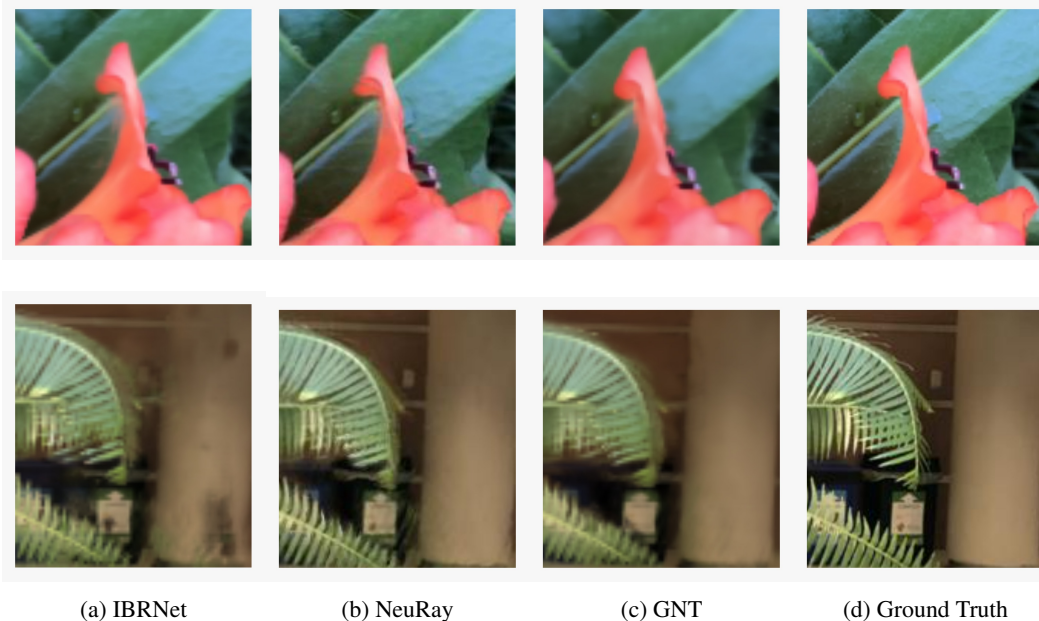


Figure 4: Qualitative results of GNT for the cross-scene rendering. On the unseen Flowers scene from the LLFF dataset (first row), GNT recovers the edges of petals more accurately. On the unseen Ferns scene from the LLFF dataset (second row), GNT is able to handle regions which are sparsely visible in the source views and hence recovers the pillar more accurately than IBRNet and visually comparable to Neuray.

(b) **Evaluation Datasets:** The evaluation datasets include the previously discussed Blender [40], LLFF datasets [62]. Please note that the LLFF scenes present in the validation set are not included in the handheld cellphone captures in the training set.

**Discussion.** We compare our method with PixelNeRF [7], MVSNerF [5], IBRNet [6], and NeuRay [46]. As seen from Tab. 3, our method outperforms SOTA by average scores of  $\sim 4\%$   $\downarrow$  in the LLFF and Blender datasets. This indicates the effectiveness of our proposed view transformer to extract generalizable scene representations. Similar to the single scene experiments, we observe weaker performance in the Blender dataset by PSNR but obtain the best perceptual metric score ( $\sim 3\%$   $\uparrow$  SSIM,  $\sim 22\%$   $\downarrow$  LPIPS). We show qualitative results in Fig. 4 and we can clearly see that GNT renders novel views with higher visual quality when compared to other methods. Specifically, as seen from the second row in Fig. 4, GNT is able to handle regions that are sparsely visible in the source views and generates images of comparable visual quality as NeuRay even with no explicit supervision for occlusion. The combination of our enhanced scene representation network - view transformer and learnable renderer - ray transformer directly contribute to these improved results.

## 5.4 Ablation Studies

To validate the significance of our architectural designs in GNT, we conduct the following ablation studies on the Drums scene from the Blender dataset.

**One-Stage Transformer:** We convert the point-wise epipolar features into one single sequence and pass it through a “one-stage transformer” network with standard dot-product self-attention layers, without considering our two-stage pipeline: view and ray aggregation.

**Epipolar Agg.  $\rightarrow$  View Agg.:** Moving to a two-stage transformer, we train a network which first aggregates features from the points along the epipolar lines followed by feature aggregation across epipolar lines on different reference views (in contrast to GNT’s first view aggregation then ray aggregation). This two-stage aggregation resembles the strategies adopted in NLF [8].

**Dot Product-Attention View Transformer:** Next we train a network that uses the standard dot-product attention in the view transformer blocks, in contrast to our proposed memory-efficient subtraction-based attention.

Table 4: Ablation study of several components in GNT on the Drums scene from the Blender dataset. The indent indicates the studied setting is added upon the upper-level ones.

Model	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	Avg $\downarrow$
One-Stage Transformer	21.80	0.862	0.152	0.072
Two-Stage Transformer				
Epipolar Agg. $\rightarrow$ View Agg.	21.57	0.863	0.153	0.073
View Agg. $\rightarrow$ Ray Agg.				
Dot Product-Attention View Transformer	26.98	<b>0.953</b>	0.089	0.034
Subtraction-Attention View Transformer				
w/ Volumetric Rendering	24.24	0.92	0.076	0.043
w/o Volumetric Rendering (Ours)	<b>29.41</b>	0.931	<b>0.085</b>	<b>0.029</b>

**w/o Volumetric Rendering:** Last but not least, we train a network to predict per-point RGB and density values from the point features output by view aggregator, and compose them using the volumetric rendering equation, instead of our attention-driven volumetric renderer.

We report the performance of the above investigation in Tab. 4. We verify that our design of two-stage transformer is superior than one-stage aggregation or an alternative two-stage pipeline [8] since our renderer strictly complies with the multi-view geometry. Compared with dot-product attention, subtraction-based attention achieves slightly higher overall scores. This also indicates the performance of GNT does not heavily rely on the choice of attention operation. What matters is bringing in attention mechanism for cross-point interaction. For practical usage, we also consider the memory efficiency in our view transformer. Our ray transformer outperforms the classic volumetric rendering implies the advantage of adopting a data-driven renderer.

### 5.5 Interpretation on Attention Maps

The use of transformers as a core element in GNT enables interpretation of the network by analyzing the attention weights. As discussed earlier, view transformer learns to find correspondence between the queried points, and neighboring views which enables it to pay attention to more “visible” views or be occlusion aware. Similarly, the ray transformer captures point-to-point interactions which enables it to model the relative importance of each point or be depth aware. We validate our hypothesis by visualization.

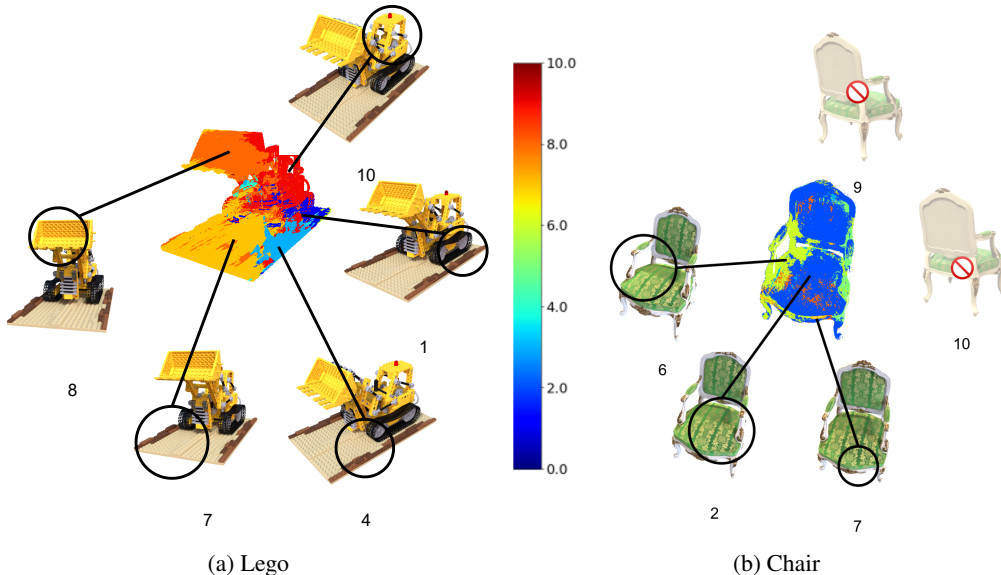


Figure 5: Visualization of view attention where each color indicates the view number which has maximum correspondence with a target pixel. GNT’s view transformer learns to map each object region in the target view to its corresponding regions in the source views which are least occluded. Transparency indicates least visible views learned via the cross-view attention.

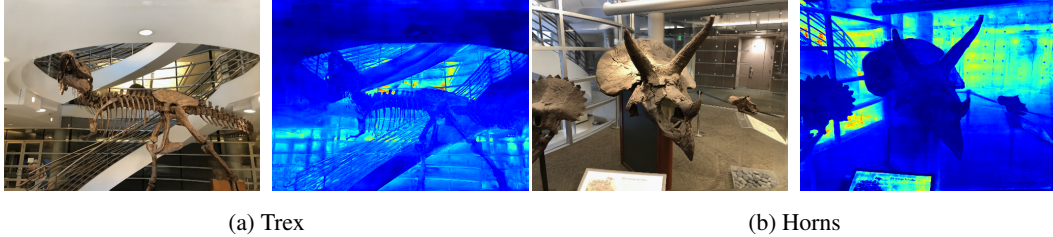


Figure 6: Visualization of ray attention where each color indicates the distance of each pixel relative to the viewing direction. GNT’s ray transformer computes point-wise aggregation weights from which the depth can be inferred. Red indicates far while blue indicates near.

**View Attention.** To visualize the view-wise attention maps learnt by our model, we use the attention matrix from Eq. 8 and collapse the channel dimension by performing mean pooling. We then identify the view-number which is assigned maximum attention with respect to each point and then compute the most repeating view number across points along a ray (by computing mode). These "view importance" values denote the view which has maximum correspondence with the target pixel’s color. Fig. 5 visualizes the source view correspondence with every pixel in the target view. Given a region in the target view, GNT attempts to pay maximum attention to a source view which is least occluded in the same region. For example: In fig. 5a, the truck’s bucket is most visible from view number 8, hence the regions corresponding to the same are orange colored, while regions towards the front of the lego are most visible from view number 7 (yellow). To further test GNT’s capacity to reason occlusion, we identify views which are farthest from the target view and corrupt the sampled source views. In Fig. 5b, every pixel in the target view (front view of chair) ignores view numbers 9, 10 (back view of the chair) which are farthest from the target view.

**Ray Attention.** To visualize the attention maps across points in a ray, we use the attention matrix from Eq. 1 and collapse the head dimension by performing mean pooling. From the derived matrix, we select a point and extract its relative importance with every other point. We then compute a depth map from these learned point-wise correspondence values by multiplying it with the marching distance of each point and sum-pooling along a ray. Fig. 6 plots the depth maps computed from the learned attention values in the ray transformer block. We can clearly see that the obtained depth maps have a physical meaning i.e pixels closer to the view directions are blue while the ones farther away are red. Therefore, with no explicit supervision, GNT learns to physically ground its attention maps.

## 6 Conclusion

We present Generalizable NeRF Transformer (GNT), a pure transformer based architecture that efficiently reconstructs NeRFs on the fly. The view transformer of GNT leverages epipolar geometry as an inductive bias for scene representation. The ray transformer renders novel views by ray marching and decoding the sequences of sampled point features using attention mechanism. Our experiments demonstrate that GNT can improve both single-scene and cross-scene training results. We also show by visualization that depth and occlusion can be inferred from attention maps. This implies that pure attention mechanism can be “universal modeling tool” for physically-grounded rendering process. Interesting future directions include relaxing the epipolar constraints to simulate more complicated light transport.

## References

- [1] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2304–2314, 2019.
- [2] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *IEEE International Conference on Computer Vision (ICCV)*, 2021.
- [3] Kai Zhang, Gernot Riegler, Noah Snaveley, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*, 2020.
- [4] Tianlong Chen, Peihao Wang, Zhiwen Fan, and Zhangyang Wang. Aug-nerf: Training stronger neural radiance fields with triple-level physically-grounded augmentations. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [5] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. Mvnerf: Fast generalizable radiance field reconstruction from multi-view stereo. In *IEEE International Conference on Computer Vision (ICCV)*, 2021.
- [6] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul P Srinivasan, Howard Zhou, Jonathan T Barron, Ricardo Martin-Brualla, Noah Snaveley, and Thomas Funkhouser. Ibrnet: Learning multi-view image-based rendering. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [7] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [8] Mohammed Suhail, Carlos Esteves, Leonid Sigal, and Ameesh Makadia. Light field neural rendering. *arXiv preprint arXiv:2112.09687*, 2021.
- [9] Nelson Max. Optical models for direct volume rendering. *IEEE Transactions on Visualization and Computer Graphics (TVCG)*, 1995.
- [10] Marc Levoy. Display of surfaces from volume data. *IEEE Computer graphics and Applications*, 8(3):29–37, 1988.
- [11] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces. *Advances in Neural Information Processing Systems*, 34, 2021.
- [12] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *arXiv preprint arXiv:2106.10689*, 2021.
- [13] Hao Chen, Bo He, Hanyu Wang, Yixuan Ren, Ser Nam Lim, and Abhinav Shrivastava. Nerv: Neural representations for videos. *Advances in Neural Information Processing Systems*, 34, 2021.
- [14] Dor Verbin, Peter Hedman, Ben Mildenhall, Todd Zickler, Jonathan T Barron, and Pratul P Srinivasan. Ref-nerf: Structured view-dependent appearance for neural radiance fields. *arXiv preprint arXiv:2112.03907*, 2021.
- [15] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems (NeurIPS)*, 2017.
- [16] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations (ICLR)*, 2020.
- [17] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International Conference on Machine Learning (ICML)*, 2021.
- [18] Yifan Jiang, Shiyu Chang, and Zhangyang Wang. Transgan: Two pure transformers can make one strong gan, and that can scale up. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.

- [19] Yikang Ding, Wentao Yuan, Qingtian Zhu, Haotian Zhang, Xiangyue Liu, Yuanjiang Wang, and Xiao Liu. Transmvsnet: Global context-aware multi-view stereo network with transformers. *arXiv preprint arXiv:2111.14600*, 2021.
- [20] Juho Lee, Yoonho Lee, Jungtaek Kim, Adam Kosior, Seungjin Choi, and Yee Whye Teh. Set transformer: A framework for attention-based permutation-invariant neural networks. In *International conference on machine learning*, pages 3744–3753. PMLR, 2019.
- [21] Ayush Tewari, Ohad Fried, Justus Thies, Vincent Sitzmann, Stephen Lombardi, Kalyan Sunkavalli, Ricardo Martin-Brualla, Tomas Simon, Jason Saragih, Matthias Nießner, et al. State of the art on neural rendering. In *Computer Graphics Forum*, volume 39, pages 701–727. Wiley Online Library, 2020.
- [22] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv:1810.04805*, 2018.
- [23] Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V Le, and Ruslan Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, 2019.
- [24] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- [25] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10012–10022, 2021.
- [26] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al. Highly accurate protein structure prediction with alphafold. *Nature*, 2021.
- [27] Chengxuan Ying, Tianle Cai, Shengjie Luo, Shuxin Zheng, Guolin Ke, Di He, Yanming Shen, and Tie-Yan Liu. Do transformers really perform badly for graph representation? *Advances in Neural Information Processing Systems*, 34, 2021.
- [28] Wenqing Zheng, Qiangqiang Guo, Hao Yang, Peihao Wang, and Zhangyang Wang. Delayed propagation transformer: A universal computation engine towards practical control in cyber-physical systems. *Advances in Neural Information Processing Systems*, 34, 2021.
- [29] Hugo Touvron, Matthieu Cord, Alexandre Sablayrolles, Gabriel Synnaeve, and Hervé Jégou. Going deeper with image transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 32–42, 2021.
- [30] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European Conference on Computer Vision (ECCV)*, 2020.
- [31] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. In *International Conference on Learning Representations (ICLR)*, 2021.
- [32] Minghang Zheng, Peng Gao, Xiaogang Wang, Hongsheng Li, and Hao Dong. End-to-end object detection with adaptive clustering transformer. In *British Machine Vision Conference (BMVC)*, 2021.
- [33] Zhiqing Sun, Shengcao Cao, Yiming Yang, and Kris Kitani. Rethinking transformer-based set prediction for object detection. In *IEEE International Conference on Computer Vision (ICCV)*, 2020.
- [34] Hanting Chen, Yunhe Wang, Tianyu Guo, Chang Xu, Yiping Deng, Zhenhua Liu, Siwei Ma, Chunjing Xu, Chao Xu, and Wen Gao. Pre-trained image processing transformer. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [35] Yuqing Wang, Zhaoliang Xu, Xinlong Wang, Chunhua Shen, Baoshan Cheng, Hao Shen, and Huaxia Xia. End-to-end video instance segmentation with transformers. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [36] Luowei Zhou, Yingbo Zhou, Jason J Corso, Richard Socher, and Caiming Xiong. End-to-end dense video captioning with masked transformer. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

- [37] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. Vivit: A video vision transformer. In *IEEE International Conference on Computer Vision (ICCV)*, 2021.
- [38] Meng-Hao Guo, Jun-Xiong Cai, Zheng-Ning Liu, Tai-Jiang Mu, Ralph R Martin, and Shi-Min Hu. Pct: Point cloud transformer. *Computational Visual Media*, 2021.
- [39] Kevin Lin, Lijuan Wang, and Zicheng Liu. End-to-end human pose and mesh reconstruction with transformers. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [40] 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 *European conference on computer vision*, pages 405–421. Springer, 2020.
- [41] 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.
- [42] Michael Oechsle, Songyou Peng, and Andreas Geiger. Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5589–5599, 2021.
- [43] Keunhong Park, Utkarsh Sinha, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Steven M. Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. *ICCV*, 2021.
- [44] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10318–10327, 2021.
- [45] Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Ricardo Martin-Brualla, and Steven M Seitz. Hypernerf: A higher-dimensional representation for topologically varying neural radiance fields. In *ACM Conference and Exhibition on Computer Graphics and Interactive Techniques in Asia (SIGGRAPH Asia)*, 2021.
- [46] Yuan Liu, Sida Peng, Lingjie Liu, Qianqian Wang, Peng Wang, Christian Theobalt, Xiaowei Zhou, and Wenping Wang. Neural rays for occlusion-aware image-based rendering. In *CVPR*, 2022.
- [47] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenotrees for real-time rendering of neural radiance fields. In *IEEE International Conference on Computer Vision (ICCV)*, 2021.
- [48] Alex Yu, Sara Fridovich-Keil, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. *arXiv preprint arXiv:2112.05131*, 2021.
- [49] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *arXiv preprint arXiv:2201.05989*, 2022.
- [50] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. In *IEEE International Conference on Computer Vision (ICCV)*, 2021.
- [51] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5459–5469, 2022.
- [52] Jeremy Reizenstein, Roman Shapovalov, Philipp Henzler, Luca Sbordone, Patrick Labatut, and David Novotny. Common objects in 3d: Large-scale learning and evaluation of real-life 3d category reconstruction. In *International Conference on Computer Vision*, 2021.
- [53] Dan Wang, Xinrui Cui, Septimiu Salcudean, and Z Jane Wang. Generalizable neural radiance fields for novel view synthesis with transformer. *arXiv preprint arXiv:2206.05375*, 2022.
- [54] Kai-En Lin, Lin Yen-Chen, Wei-Sheng Lai, Tsung-Yi Lin, Yi-Chang Shih, and Ravi Ramamoorthi. Vision transformer for nerf-based view synthesis from a single input image. *arXiv preprint arXiv:2207.05736*, 2022.

- [55] Mehdi SM Sajjadi, Henning Meyer, Etienne Pot, Urs Bergmann, Klaus Greff, Noha Radwan, Suhani Vora, Mario Lučić, Daniel Duckworth, Alexey Dosovitskiy, et al. Scene representation transformer: Geometry-free novel view synthesis through set-latent scene representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6229–6238, 2022.
- [56] Mehdi SM Sajjadi, Daniel Duckworth, Aravindh Mahendran, Sjoerd van Steenkiste, Filip Pavetić, Mario Lučić, Leonidas J Guibas, Klaus Greff, and Thomas Kipf. Object scene representation transformer. *arXiv preprint arXiv:2206.06922*, 2022.
- [57] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015.
- [58] Hengshuang Zhao, Li Jiang, Jiaya Jia, Philip HS Torr, and Vladlen Koltun. Point transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 16259–16268, 2021.
- [59] Zhiwen Fan, Tianlong Chen, Peihao Wang, and Zhangyang Wang. Cadtransformer: Panoptic symbol spotting transformer for cad drawings. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10986–10996, 2022.
- [60] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. In *Proceedings of the European conference on computer vision (ECCV)*, pages 767–783, 2018.
- [61] Zhenpei Yang, Zhile Ren, Qi Shan, and Qixing Huang. Mvs2d: Efficient multi-view stereo via attention-driven 2d convolutions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8574–8584, 2022.
- [62] 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.
- [63] Suttisak Wizadwongsa, Pakkapon Phongthawee, Jiraphon Yenphraphai, and Supasorn Suwanakorn. Nex: Real-time view synthesis with neural basis expansion. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [64] Ellen D Zhong, Tristan Bepler, Bonnie Berger, and Joseph H Davis. Cryodrgn: reconstruction of heterogeneous cryo-em structures using neural networks. *Nature Methods*, 18(2):176–185, 2021.
- [65] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004.
- [66] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018.
- [67] Kai-En Lin, Lin Yen-Chen, Wei-Sheng Lai, Tsung-Yi Lin, Yi-Chang Shih, and Ravi Ramamoorthi. Vision transformer for nerf-based view synthesis from a single input image. In *arxiv*, 2022.
- [68] Laura Downs, Anthony Francis, Nate Koenig, Brandon Kinman, Ryan Hickman, Krista Reymann, Thomas B McHugh, and Vincent Vanhoucke. Google scanned objects: A high-quality dataset of 3d scanned household items. *arXiv e-prints*, pages arXiv–2204, 2022.
- [69] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: Learning view synthesis using multiplane images. *arXiv preprint arXiv:1805.09817*, 2018.
- [70] John Flynn, Michael Broxton, Paul Debevec, Matthew DuVall, Graham Fyffe, Ryan Overbeck, Noah Snavely, and Richard Tucker. Deepview: View synthesis with learned gradient descent. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2367–2376, 2019.



## A Implementation

---

**Algorithm 1** Cross View Attention: PyTorch-like Pseudocode

---

$X_o \rightarrow$  coordinate aligned features( $N_{\text{rays}}, N_{\text{pts}}, D$ )  
 $X_j \rightarrow$  epipolar view features( $N_{\text{rays}}, N_{\text{pts}}, N_{\text{views}}, D$ )  
 $\Delta d \rightarrow$  relative directions of source views wrt target views( $N_{\text{rays}}, N_{\text{pts}}, N_{\text{views}}, 3$ )  
 $f_Q, f_K, f_V, f_P, f_A, f_O \rightarrow$  functions that parameterize MLP layers

$Q = f_Q(X_o)$   
 $K = f_K(X_j)$   
 $V = f_V(X_j)$

$P = f_P(\Delta d)$   
 $A = K - Q[:, :, \text{None}, :] + P$   
 $A = \text{softmax}(A, \text{dim} = -2)$

$O = ((V + P) \cdot A). \text{sum}(\text{dim} = 2)$   
 $O = f_O(O)$

---

---

**Algorithm 2** Ray Attention: PyTorch-like Pseudocode

---

$X_o \rightarrow$  coordinate aligned features( $N_{\text{rays}}, N_{\text{pts}}, D$ )  
 $x \rightarrow$  point coordinates (after position encoding)( $N_{\text{rays}}, N_{\text{pts}}, D$ )  
 $d \rightarrow$  target view direction (after position encoding)( $N_{\text{rays}}, N_{\text{pts}}, D$ )  
 $f_Q, f_K, f_V, f_P, f_A, f_O \rightarrow$  functions that parameterize MLP layers

$X_o = f_P(\text{concat}(X_o, d, x))$   
 $Q = f_Q(X_o)$   
 $K = f_K(X_o)$   
 $V = f_V(X_o)$

$A = Q \cdot K^T / \sqrt{D}$   
 $A = \text{softmax}(A, \text{dim} = -1)$

$O = (V \cdot A)$   
 $O = f_O(O)$

---

We provide a simple and efficient pytorch pseudo-code to implement the attention operations in the view, ray transformer blocks in Alg. 1, 2 respectively. We do not indicate the feedforward and layer normalization operations for simplicity. As seen in Alg. 3, we reuse the epipolar view features  $X_j$  to derive keys, and values across view transformer blocks. Therefore, one could further improve efficiency by computing them only once while also sharing the network weights across view transformer blocks or simply put  $f_{\text{view } i}(\cdot)$  represents the same function across different values of  $i$ . This can be considered analogous to an unfolded recurrent neural network that updates itself iteratively but using the same weights.

## B Extensions

### B.1 Auto-Regressive Decoding

The final rendered color is obtained by mean-pooling the outputs from the ray transformer block, and mapping the pooled feature vector to RGB via an MLP layer. It can be understood that the target pixel’s color is strongly dependent on the closest point from the ray origin and weakly related to the farthest. Revisiting Eq. 9, volumetric rendering attempts to compose point-wise color depending on the other points in a far to near fashion. Motivated by this, we propose an auto-regressive decoder to better simulate the rendering process. Transformers have shown great success in auto-regressive

---

**Algorithm 3** GNT: PyTorch-like Pseudocode
 

---

$X_j \rightarrow$  epipolar view features ( $N_{\text{rays}}, N_{\text{pts}}, N_{\text{views}}, D$ )  
 $x \rightarrow$  point coordinates (after position encoding) ( $N_{\text{rays}}, N_{\text{pts}}, D$ )  
 $d \rightarrow$  target view direction (after position encoding) ( $N_{\text{rays}}, N_{\text{pts}}, D$ )  
 $\Delta d \rightarrow$  relative directions of source views wrt target views ( $N_{\text{rays}}, N_{\text{pts}}, N_{\text{views}}, 3$ )  
 $f_{\text{view } i}, f_{\text{ray } i} \rightarrow$  functions that parameterize view transforms, ray transformers at layer  $i$  respectively  
 $f_{\text{rgb}} \rightarrow$  functions that parameterize MLP layers

```

i = 0
X_o = max_{j=1}^{N_{\text{views}}} (X_j)
while i < N_{\text{layers}} do
  X_o = f_{\text{view } i}(X_o, X_j, \Delta d)
  X_o = f_{\text{ray } i}(X_o, d, x)
  i = i + 1
end while
O = Norm(X_o)
RGB = f_{\text{rgb}}(\text{mean}_{i=1}^{N_{\text{points}}} (X_o))
  
```

---

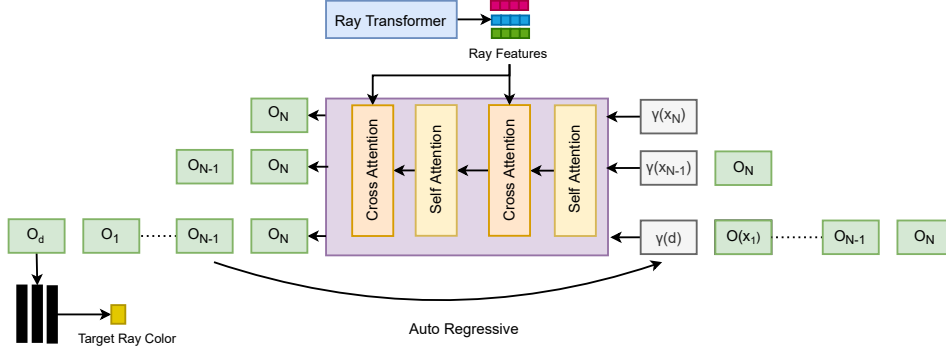


Figure 7: Architecture of Auto-Regressive Decoder and Sampling strategy in a far to near fashion.

decoding, more specifically in NLP [15]. We borrow a similar strategy and replace the simpler MLP based color prediction with a series of transformer blocks - with self, cross attention layers.

In the first pass, the decoder is queried with positional encoding of the farthest point ( $\gamma(x_N)$ ) to generate an output feature representation of the same. In the next step, the output token is concatenated with the second farthest point ( $\gamma(x_{N-1})$ ) to query the decoder. This process repeats until all the points in the ray are queried in a far to near fashion. In the final pass, the encoded view direction ( $\gamma(d)$ ) is concatenated with the per-point output features in the previous passes to query the decoder and the output token corresponding to the view direction is extracted. The extracted token is mapped to RGB via an MLP layer. This entire process is summarized in Fig. 7. The auto-regressive procedure closely resembles the volumetric rendering equation which iteratively blends and overrides the previous color when marching along a ray from far to near.

Transformer-based decoders used in language iteratively predict output tokens only during inference, i.e. they are trained in a non-autoregressive fashion due to the availability of ground truth output tokens in each step. Transferring the same to neural rendering is not possible, as we do not have access to ground truth color for each point sampled along the ray. Hence, we require a loop to auto-regressively decode features even during training. This reduces the computational efficiency of the proposed strategy especially as the number of points sampled along the ray increases. Therefore, we introduce a caching mechanism to store the per-layer outputs of the previous tokens and only compute the attention of a new token in the current pass. This does not overcome the iterative loop during each forward pass but avoid redundant computations, which helps improve the decoding speed drastically when compared to the naive strategy. Due to computational constraints, we are only able to train GNT + AutoReg with much fewer rays sampled per iteration (500) when compared to

Table 5: Comparison of auto regressive GNT against SOTA methods for single scene rendering on the LLFF dataset.

Models	Orchids				Horns			
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	Avg $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	Avg $\downarrow$
LLFF	18.52	0.588	0.141	0.313	24.70	0.840	0.193	0.064
NeRF	20.36	0.641	0.321	0.121	27.45	0.828	0.268	0.058
NeX	20.42	0.765	0.242	0.102	28.46	0.934	0.173	0.040
NLF	21.05	<b>0.807</b>	0.173	0.084	<b>29.78</b>	<b>0.957</b>	0.121	0.030
GNT	<b>21.20</b>	0.760	<b>0.161</b>	<b>0.084</b>	29.41	0.931	0.085	<b>0.029</b>
GNT + AutoReg	21.05	0.736	0.181	0.090	28.20	0.908	0.114	0.037
GNT + Fine	21.16	0.756	0.164	0.085	28.96	0.928	<b>0.083</b>	0.030

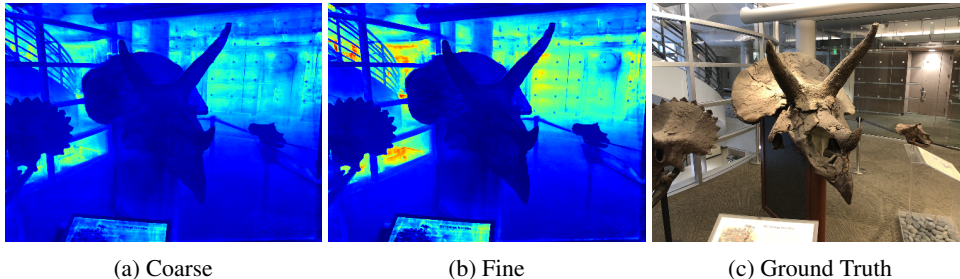


Figure 8: Visualization of ray attention extracted during coarse, fine sampling where each color indicates the distance of each pixel relative to the viewing direction. The sampled fine points inferred from the learn attention values help GNT capture more fine-grained details. Red indicates far while blue indicates near.

other methods as discussed in Sec. 5.2. Tab. 5 discusses single scene optimization results on the LLFF dataset and we can clearly see that the GNT + AutoReg improves the overall performance when compared to existing baselines. However, GNT with the auto-regressive decoder does not improve performance when compared to our own method without the decoder. This could be because of the fewer number of rays sampled and we expect our results would improve when scaled to similar settings. Nevertheless, this shows that the learnable decoder predicts per-point RGB features effectively without any supervision from the volumetric rendering equation.

## B.2 Attention Guided Coarse-to-Fine Sampling

GNT’s ray transformer learns point-to-point correspondences which helps model visibility and occlusion or more formally point-wise density  $\sigma$ . Motivated by this hypothesis, we estimate the depth maps from the extracted attention maps and qualitatively analyze the same in Sec. 5.5. Therefore, we could conclude that the learned point-wise importance values can be considered equivalent to point-wise density or  $\sigma$ . To further test our claim, we attempt to use the learned point-wise correspondence values to sample “fine” points, which are then queried to GNT to render a higher-quality image. Due to the set-like property of attention, we directly query the fine points to the same network without using a separately trained “fine” network unlike other NeRF methods [40, 2, 6, 46]. Please note that we follow the same training strategy from Sec. 5.1 and only sample coarse-fine points during evaluation. Tab. 5 compares “GNT + Fine” against other methods, and we can clearly see that it outperforms other SOTA methods in complex scenes like Orchids and on-par in other scenes like Horns. However, it does not outperform our own method without fine sampling. We attribute this performance gap to the lack of training with the coarse-fine sampling strategy and expect our results to improve further. In Fig. 8, we visualize the estimated depth values obtained from the learned attention maps during both coarse, and fine stages. We can clearly see that the fine depth map is able to better estimate differences between nearby pixels which results in a higher resolution output.