

Lightplane: Highly-Scalable Components for Neural 3D Fields

Ang Cao^{1,2}, Justin Johnson², Andrea Vedaldi¹, and David Novotny¹

¹ Meta AI

² University of Michigan

Abstract. Contemporary 3D research, particularly in reconstruction and generation, heavily relies on 2D images for inputs or supervision. However, current designs for these 2D-3D mapping are memory-intensive, posing a significant bottleneck for existing methods and hindering new applications. In response, we propose a pair of highly scalable components for 3D neural fields: *Lightplane Renderer* and *Splatter*, which significantly reduce memory usage in 2D-3D mapping. These innovations enable the processing of vastly more and higher resolution images with small memory and computational costs. We demonstrate their utility in various applications, from benefiting single-scene optimization with image-level losses to realizing a versatile pipeline for dramatically scaling 3D reconstruction and generation. Code: <https://github.com/facebookresearch/lightplane>.

1 Introduction

Recent advancements in neural rendering and generative modeling have propelled significant strides in 3D reconstruction and generation. However, in most

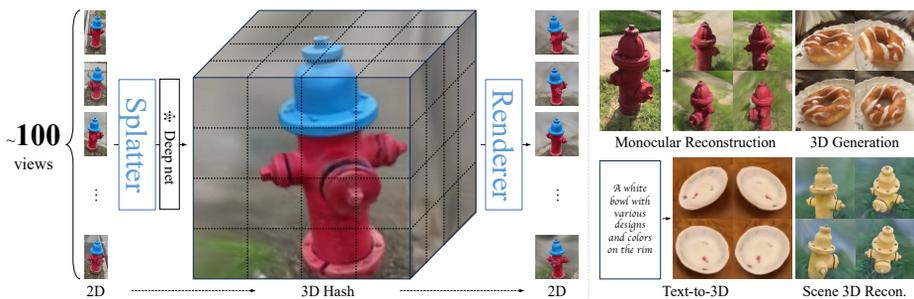


Fig. 1: We introduce the *Lightplane Renderer* and *Splatter*, a pair of highly-scalable components for neural 3D fields (left). They address the key memory bottleneck of 2D-3D mapping (*i.e.* rendering and lifting), and reduce memory usage by up to four orders of magnitude, which dramatically increases the number of images that can be processed. We showcase how they can boost various 3D applications (right).

cases, these methods are not exclusively 3D; instead, they heavily rely on 2D images as inputs or for supervision, which demands information mapping between 2D and 3D spaces. For instance, Neural Radiance Fields (NeRFs) [54] use a photometric loss on 2D images rendered from 3D, bypassing direct 3D supervision. Similarly, various novel view synthesis and generation methods [11, 15, 98] employ 2D images as inputs and lift them into 3D space for further processing. This mapping between 2D and 3D is critical for current 3D research, attributed to the scarcity of 3D training materials for developing versatile 3D models from scratch, and the relative ease of acquiring 2D images on a large scale.

Despite its crucial role and widespread use, the process of 2D-3D mapping incurs a high computational cost, especially in *neural 3D fields* with volumetric rendering, which underpins many of the most powerful 3D representations. These fields are defined by continuous functions that assign values, such as density or color, to *any* point in 3D space, regardless of the presence of a physical surface. Therefore, they are powerful and flexible, preventing initialization in point rendering [38] or topology constraints for meshes [50]. The primary challenge lies in executing operations across numerous 3D points that span an entire volume. While these operations can be relatively simple (*e.g.*, evaluating a small multilayer perceptron (MLP) at each point, or extracting features from 2D input feature maps), performing them in a differentiable manner is *extremely* memory intensive as all intermediate values must be kept in memory for backpropagation.

While the speed of NeRFs has been improved in [12, 22, 56], the issue of high memory consumption has seldom been studied. This significant memory demand hampers scalability of 2D-3D communication, presenting a crucial bottleneck for many existing 3D models and a formidable barrier for potential new applications. For example, the memory requirements to render even a single low-resolution image of a neural 3D field can be prohibitive enough to prevent the application of image-level losses such as LPIPS [102] or SDS [61]. Omitting such losses leads to a massive performance loss, as *e.g.* demonstrated by the state-of-the-art Large Reconstruction Model [3, 31]. Additionally, memory-inefficiencies limit the number of input images and the resolution of the 3D representation, preventing advancing from few-view novel view synthesis models [88, 98] to large-scale amortized 3D reconstruction models with many conditioning images.

In this paper, we propose two highly scalable components for neural 3D fields: *Lightplane Renderer* and *Splatter*. These innovations enable 2D-3D mapping with four orders of magnitude less memory consumption while maintaining comparable speed. *Renderer* renders 2D images of 3D models by means of the standard emission-absorption equations popularized by NeRF [54]. Conversely, *Splatter* lifts 2D information to 3D by splatting it onto the 3D representation, allowing further processing with neural nets. Both components are based on a hybrid 3D representation that combines ‘hashed’ 3D representations such as voxel grids and triplanes with MLPs. We use these representations as they are fast, relatively memory efficient, and familiar to practitioners, while components could be easily extended to other hashed representations as well.

As aforementioned, storing intermediate values at each 3D point for back-propagation causes tremendous memory usage. We solve it by creatively re-configuring inner computations and fusing operations over casted rays instead of 3D points. Specifically, *Lightplane Renderer* sequentially calculates features (e.g., colors) and densities of points along the ray, updating rendered pixels and transmittance on-the-fly without storing intermediate tensors. This design significantly saves memory at the cost of a challenging backpropagation, which we solve by efficiently recomputing forward activations as needed. Note that the latter is different from the standard “checkpointing” trick, whose adoption here would be of little help. This is because checkpointing still entails caching many intermediate ray-point values as we march along each ray.

Lightplane Splatter builds on similar ideas with an innovative design, where splatted features are stored directly into the hash structure underpinning the 3D model, without emitting one value per 3D point. Besides voxel grids which are usually used for lifting, *Splatter* could be easily extended to other 3D hash structures. We implement these components in Triton [81], a GPU programming language that is efficient, portable, and relatively easy to modify. We will release the code as an open-source package upon publication.

Like convolution or attention, our components are designed as building blocks to boost a variety of 3D models and applications. Empowered by the *Lightplane*, we devise a pipeline taking up to input 100 images, significantly scaling the communication between 2D and 3D. We extensively evaluate on the CO3Dv2 dataset, reporting significant performance improvements in color and geometry accuracy for 3D reconstruction, and better 3D generation measured by FID/KID. Finally, we boost performance of the state-of-the-art Large Reconstruction Model [31].

2 Related Work

3D reconstruction using neural 3D fields. Traditional 3D reconstruction models represented shapes as meshes [26, 86], point clouds [20, 95], or voxel grids [16, 25]. With the introduction of NeRF [54], however, the focus has shifted to *implicit* 3D representations, often utilizing MLPs to represent occupancy and radiance functions defined on a 3D domain. NeRF has been refined in many ways [4, 5, 85, 101], including replacing the opacity function with a signed-distance field to improve the reconstruction of surfaces [46, 70, 87, 91, 96, 99].

Storing an entire scene in a single MLP, however, means evaluating a complex function anew at every 3D point, which is very expensive from both time and memory usage. Many authors have proposed to represent radiance fields with smaller, more local components to improve speed, including using point clouds [94], tetrahedral meshes [40] or, more often, voxel grids [36, 51, 64, 77, 97]. Voxel grids could be further replaced by more compact structures like low-rank tensor decompositions [13], triplanes [9], hashing [57], and their combination [65].

Unlike the above methods focusing on speed, *Lightplane* significantly reduces memory demands for neural 3D fields. Note that our method targets neural 3D fields with volumetric rendering, while point-based rendering like 3DGS [38] are

not in this scope, since they don’t model every 3D point in the space and rely on rasterization instead of volumetric rendering. While 3DGS exhibits fast convergence speed, importantly, it has been shown to give lower accuracy (measured in PSNR) in both the single-scene overfitting case [6], and the few-view reconstruction case [82]. For optimal performance in single-scene, they require careful surface initialization whereas NeRFs converge from a random initialization.

Amortized 3D reconstruction. Amortized (Generalizable) 3D reconstruction utilizing implicit shape representations was initially approached in [30, 58, 66, 83, 89, 98] by warping/pooling features from source views to a target to estimate the color of the underlying scene surfaces. [71, 92] introduces latent transformer tokens to support the reconstruction. Generalizable triplanes [31, 32, 42], ground-planes [73], and voxel grids [34] were also explored.

A common downside of these methods is their memory consumption which limits them all to a *few-view* setting with up to 10 source views. They either are trained on a category-specific dataset or learn to interpolate between input views with unsatisfactory geometry and 3D consistency. Owing to its memory efficiency, *Lightplane* allows more than 100 input source views. We leverage the latter to train a large-scale 3D model yielding more accurate reconstructions.

Image-supervised 3D generators. With the advent of Generative Adversarial Networks [27] (GAN), many methods attempted to learn generative models of 3D shapes given large uncensored image datasets. PlatonicGAN [29], HoloGAN [59] and PrGAN [23] learned to generate voxel grids whose renders were indistinguishable from real object views according to an image-based deep discriminator. The same task was later tackled with Neural Radiance Fields [9, 28, 60, 72, 74], and with meshes [24, 93]. The success of 2D generative diffusion models [19] led to image-supervised models such as HoloDiffusion [37], Forward Diffusion [80], and PC² [53], which directly model the distribution of 3D voxel grids, implicit fields and point clouds respectively. Similarly, RenderDiffusion [1] and ViewsetDiffusion [79] learn a 2D image denoiser by means of a 3D deep reconstructor. GenVS [11] and HoloFusion [35] proposed 3D generators with 2D diffusion rendering post-processors. We demonstrate that *Lightplane* brings a strong performance boost to ViewsetDiffusion and generates realistic 3D scenes.

3 Method

We introduce the *Lightplane Renderer* and *Splatter*, which facilitate the mapping of information between 2D and 3D spaces in a differentiable manner, significantly reducing memory usage in the process. We first discuss the memory bottlenecks of existing methods that are used for rendering and lifting images into 3D structures (Sec. 3.1). Then we define the hashed 3D representations (Sec. 3.2) used in our framework and functionality (Sec. 3.3) of the proposed components. Lastly, we discuss their implementations (Sec. 3.4).

3.1 Preliminary

2D-3D Mapping. Mapping between 2D images and 3D models is a major practical bottleneck of many algorithms (Sec. 1), particularly when using powerful implicit 3D representations such as neural 3D fields. The memory bottleneck comprises a large number of 3D points from rendering rays and their intermediate features, which are cached in GPU memory for the ensuing backpropagation.

More specifically, for rendering (*3D to 2D mapping*), an *entire ray* of 3D points contributes to the color of a single pixel in the rendered image. With M pixels and R points per ray, $M \times R$ implicit representation evaluations are required to get 3D points’ colors and opacities. All these intermediate results, including outputs of all MLP layers for every 3D point, are stored in memory for backpropagation, leading to huge memory usage. Using a tiny MLP with $L=6$ layers and $K=64$ hidden units, $M \times R \times L \times K$ memory is required to *just* store the MLP outputs, which totals 12 GB for a 256^2 image with $R=128$ points per ray.

Similarly, to lift N input features to 3D (*2D to 3D mapping*), popular models like PixelNeRF [98] and GeNVS [11] project each 3D point to N input views individually, and average N sampled feature vectors as the point feature. Even without considering any MLPs, $N \times |\mathcal{M}|$ memory is used, where $|\mathcal{M}|$ is the size of 3D structure \mathcal{M} . When \mathcal{M} is a 128^3 voxel grid with 64-dimensional features, $|\mathcal{M}|$ takes 512 MB in FP32, leading to 5 GB of memory with just 10 input views.

Moreover, the aforementioned lifting requires 3D positions for projection and cannot be easily generalized to other compact representations like triplanes, since cells in such “hashed” feature maps (*e.g.* 2D position on feature planes for triplanes) don’t have clearly-defined 3D positions. Hence, directly lifting multi-view features to triplanes for further processing is still an open problem.

The memory bottleneck impacts several aspects. For mapping from 3D to 2D (*i.e.* rendering), methods like NeRF [54] and PixelNeRF [98] are limited to a few low-resolution images per training iteration (even using 40GB GPUs) or to sub-sample rendered pixels, which prohibits image-level losses such as LPIPS [102] and SDS [61]. For mapping from 2D to 3D, memory demands limit input view numbers and 3D representation sizes. The huge memory usage not only occupies resources that could otherwise enhance model sizes and capacities but also restricts model training and inference on devices with limited memory availability.

Neural 3D fields. Let $\mathbf{x} \in \mathbb{R}^3$ denote a 3D point, a *neural 3D field* is a volumetric function f that maps each point \mathbf{x} to a vector $f(\mathbf{x}) \in \mathbb{R}^C$. NeRF [54]

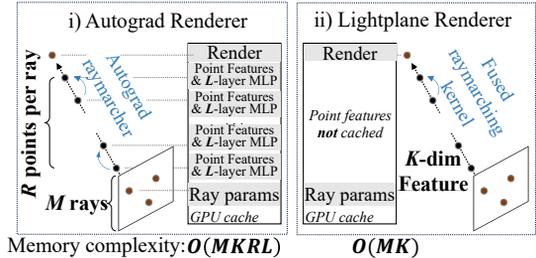


Fig. 2: Memory usage of our *Lightplane Renderer* vs. a standard autograd NeRF renderer.

represented such functions using a single MLP. While this is simple, the MLP must represent *whole* 3D objects and hence must be large and costly to evaluate.

Several approaches are proposed to solve this problem, by decomposing information into local buckets, accessing which is more efficient than evaluating the global MLP. Most famously, [56] utilizes hash tables, but other representations such as voxel grids [78] and various low-rank decompositions such as triplanes [9], TensorRF [14] and HexPlane [7, 10] also follow this pattern.

3.2 Hybrid representation with 3D hash structure

Following the idea in Sec. 3.1, we use a hybrid representation for neural 3D fields f , and decompose $f = g \circ h$, where $h : \mathbb{R}^3 \rightarrow \mathbb{R}^K$ is a hashing scheme (sampling operation) for 3D hash structure θ , and $g : \mathbb{R}^K \rightarrow \mathbb{R}^C$ is a tiny MLP, which takes features from hashing as inputs and outputs the final values. In this paper, we generalize the concept of 3D hash structures to structures like voxel grids [78], triplanes [9], HexPlane [7, 10] and actual hash table [56], as obtaining information from these structures only requires accessing and processing the small amount of information stored in a particular bucket. The associated hashing scheme h typically samples 3D point features from hash structure θ via interpolation, which is highly efficient. In practice, we operationalize θ with voxel grids and triplanes as they are easy to process by neural networks, although other structures with a differentiable hashing scheme could be easily supported.

In more detail, in the voxel-based representation, θ is a $H \times W \times D \times K$ tensor and h is the tri-linear interpolation on θ given position \mathbf{x} . In the triplane representation, θ is a list of three tensors of dimensions $H \times W \times K$, $W \times D \times K$, and $D \times H \times K$. Then, $h(x, y, z)$ is obtained by bilinear interpolation of each plane at (x, y) , (y, z) , (z, x) , respectively, followed by summing the resulting three feature vectors. Again, this design could be easily generalized to other hashed 3D structures θ and their corresponding hashing scheme (sampling operation) h .

3.3 Rendering and splatting

We now detail *Lightplane Renderer* and *Splatter*, two components using hybrid 3D representations with 3D hash structures. They are mutually dual as one maps 3D information to 2D via rendering, and the other maps 2D images to 3D.

Renderer. *Renderer* outputs pixel features \mathbf{v} (e.g. colors, depths) in a differentiable way from a hybrid representation $f = g \circ h$, given M rays $\{\mathbf{r}_i\}_{i=1}^M$ and $R+1$ points per ray. We make its high-level design consistent with existing hybrid representations [7, 9, 56, 78] as they have proven to be powerful, while re-designing the implementation in Sec. 3.4 to achieve significant memory savings.

Following volumetric rendering of NeRF [54], *Renderer* uses a generalized Emission-Absorption (EA) model and calculates transmittance T_{ij} , which is the probability that a photon emitting at \mathbf{x}_{ij} (j -th sampling points on the i -th ray)

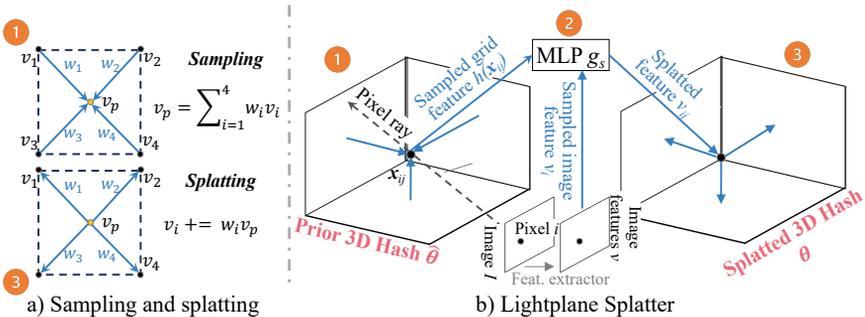


Fig. 3: Lightplane Splatter. (a) On a hash grid with vertex features \mathbf{v}_i : *sampling* obtains point features \mathbf{v}_p by interpolating vertex features weighted by inverse distance; *splatting* updates vertex features by accumulating point feature to vertex using the same weights. (b) *Splatter* involves three steps. For each 3D point along the ray, *Splatter* samples its features from prior 3D hash $\hat{\theta}$ (1), calculates features to be splatted using MLP (2), and splats them to zero-initialized θ (3).

reaches the sensor. Accordingly, the rendered feature \mathbf{v}_i of ray \mathbf{r}_i is:

$$\mathbf{v}_i = \sum_{j=1}^R (T_{i,j-1} - T_{ij}) f_v(\mathbf{x}_{ij}). \quad (1)$$

where $f_v(\mathbf{x}_{ij})$ is the feature (e.g. color) of the 3D point \mathbf{x}_{ij} , obtained from the hybrid representation f_v ; $T_{ij} = \exp(-\sum_{n=0}^j \Delta \cdot \sigma(\mathbf{x}_{in}))$, Δ is the distance between two sampled points, and $\sigma(\mathbf{x}_{in})$ is the opacity of the n -th sampled point; $(T_{i,j-1} - T_{ij}) \in [0, 1]$ is the visibility of the point \mathbf{x}_{ij} . Given a 3D point, *Renderer* samples its feature from the 3D representation and feeds the feature to an MLP g_σ to calculate the opacity. $f_v(\mathbf{x}_{ij})$ is calculated by another MLP g_v taking the sampled feature and view directions as inputs.

Splatter. Opposite to *Renderer*, *Splatter* maps input view features to 3D hash structures. Existing works like [11,35,37,79] achieve this by looping over all points *inside voxel grids* and *pulling* information from input features. They project 3D points to input views, interpolate fields of 2D image features, compute and store a feature vector for each 3D sample. Such operations are inherently memory-intensive and cannot be easily generalized to other 3D hash structures Sec. 3.1.

Instead of looping over 3D points and pulling information from inputs, we make *Splatter* loop over *input pixels/rays* and directly *push* information to 3D structures. This makes *Splatter* a reversion of *Renderer*, being able to easily extend to other 3D structures and enjoy similar memory optimization designs.

Given M input pixels, *Splatter* expands each pixel into a ray \mathbf{r}_i with $R + 1$ equispaced 3D points \mathbf{x}_{ij} , with points along the ray inheriting the pixel’s features \mathbf{v}_i . 3D points’ features \mathbf{v}_{ij} are splatted back to zero-initialized 3D structures θ , which operation is inverse to the sampling operation $h(\mathbf{x})$ used in rendering. This is done by accumulating \mathbf{v}_{ij} to hash cells that contain \mathbf{x}_{ij} , which accumulation is weighted by splatting weights. After accumulating over all M rays, each hash cell is normalized by the sum of all splatting weights landing in the cell. The

splatting weights are the same as the sampling weights used in rendering. For voxel grids, a hash cell is a voxel, and splatting weights are the normalized inverse distance between the 3D point and eight voxel vertices. It can be easily extended to other hash structures. We illustrate this splatting operation in Figure 3(b).

This naïve version of *Splatter* works well for voxel grids, but fails to work on triplanes and potentially other hashed representations. We hypothesize it is due to 3D position information being destroyed when reducing from 3D space to 2D planes, and accumulated features are unaware of the spatial structure of the 3D points. To address this, we propose to use an MLP g_s to predict a modified feature vector \mathbf{v}_{ij} from the input vectors \mathbf{v}_i , interpolated prior shape encoding $h_{\hat{\theta}}(\mathbf{x}_{ij})$, and the positional encoding $\text{direnc}(\mathbf{r}_i)$ of ray direction \mathbf{r}_{ij} . For each sample \mathbf{x}_{ij} , the splatted feature $\tilde{\mathbf{v}}_{ij}$ is

$$\tilde{\mathbf{v}}_{ij} = g_s(\mathbf{v}_i, h_{\hat{\theta}}(\mathbf{x}_{ij}), \text{direnc}(\mathbf{r}_{ij})) \quad (2)$$

$\hat{\theta}$ is another hashed 3D representation, where prior shape encoding of 3D point \mathbf{x}_{ij} could be obtained by hashing operation $h_{\hat{\theta}}(\cdot)$. This MLP allows points along the same ray to have different spatial-aware features and thus preserves the spatial structure of the 3D points. This design also allows us to iteratively refine 3D representations θ based on previous representations $\hat{\theta}$ and input features.

3.4 Memory-efficient Implementation

We discuss the practical implementations of *Lightplane Renderer* and *Splatter*, which are designed to be memory-efficient and scalable.

Fusing operations along the ray. As analyzed, current rendering and lifting operations for neural 3D fields are memory intensive, as they treat 3D points as basic entities and store intermediate results for each point. Alternatively, we treat rays as basic entities and fuse operations in a single GPU kernel, where each kernel instance is responsible for a single ray. This allows us to only store the rendered features and accumulated transmittance of the ray.

As Eq. 1, a *Renderer* kernel sequentially samples 3D points’ features, calculates features and opacities via MLPs and updates the rendered results and accumulated transmittance of the ray. These processes are integrated into a single kernel, obviating the need for storing any other intermediate results. For the example in Sec. 3.1, memory usage is significantly reduced from $O(MKRL)$ to $O(MK)$, decreasing from 12 GB to 2 KB for an image of size 256^2 with $R = 128$ samples per ray in FP32. This is less than 0.02% of the memory required by the naïve implementation. Since *Splatter* is designed to process rays emanating from input pixels as well (Sec. 3.3), it benefits from the same optimization practice.

Recalculation for backpropagation. Saving *no* intermediate results during forward propagation significantly decreases memory usage, while these tensors are essential for backpropagation. To solve it, we recompute the intermediate results during backpropagation for gradient calculation. Speed-wise, recalculating the MLP in the forward direction increases the total number of floating-point

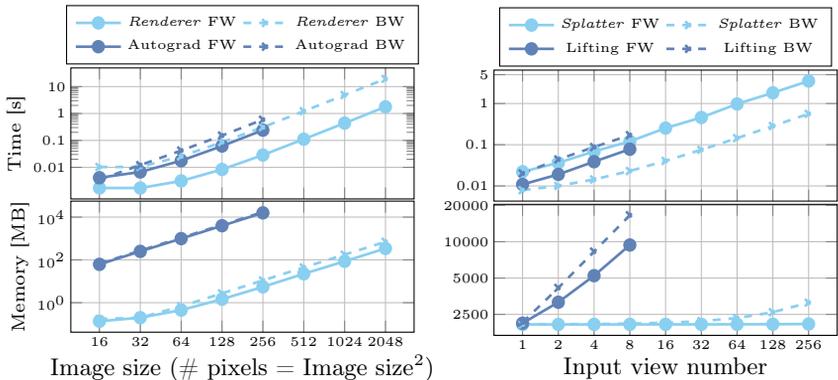


Fig. 4: *Lightplane* memory & speed benchmark showing the forward (FW and backward (BW) passes of *Lightplane Renderer* (left) and *Splatter* (right), compared to the *Autograd* renderer and *lifter* from [11, 98]. *Lightplane* exhibits up to 4 orders of magnitude lower memory consumption at comparable speed. All axes are log-scaled.

operations by less than 50% compared to the naïve implementation. But this cost only occurs during backpropagation, and leads to massive memory savings.

Leveraging GPU memory hierarchy for speed. The speed could be further optimized by exploiting the hierarchical architecture of GPU memory. By fusing operations in a single GPU kernel, we enhance the utilization of GPU’s on-chip SRAM, and prevent massive access to GPU’s high bandwidth memory (HBM). Given that HBM access speeds are substantially slower compared to on-chip SRAM, and performance bottlenecks often stem from HBM access during tensor read/write operations, our kernel maintains a competitive speed even while recalculating intermediate results for backpropagation. We encourage readers to refer to flash-attention [17] for details of GPU memory hierarchy.

Emission-absorption backpropagation. *Renderer* and *Splatter* are dual to each other not only in functionality but also in their high-level implementation. The backpropagation process of *Splatter* mirrors the forward pass of *Renderer*, as it samples 3D point gradients from the representation’s gradient field and aggregates them along the ray to form the input pixel’s gradients. Conversely, *Renderer*’s backward process is also similar to *Splatter*’s forward pass.

Notably, the backpropagation of *Renderer* is more complicated as the visibility of 3D points is affected by the transmittance of previous points in the emission-absorption model. During forward pass, we sequentially calculate 3D points’ visibility and implement the rendering equation Eq. (1) by summing in order $j = 1, 2, \dots$, as it is easy to obtain visibility T_j from T_{j-1} (we omit ray index for simplicity). For backward pass, on a ray \mathbf{r} , we derive the vector-Jacobian

product (*i.e.*, the quantity computed during backpropagation) of *Renderer*:

$$\mathbf{p}^\top \frac{d\mathbf{v}}{df_\sigma(\mathbf{x}_q)} = -\Delta \frac{d\sigma(\mathbf{x}_q)}{df_\sigma(\mathbf{x}_q)} \left(\sum_{j=q+1}^R (T_{j-1} - T_j) \mathbf{a}_j - T_q \mathbf{a}_q \right), \quad (3)$$

where $\mathbf{a}_j = \mathbf{p}^\top f_v(\mathbf{x}_j)$ and \mathbf{p} is the gradient vector that needs backpropagating.

To backpropagate through *Renderer* efficiently, we compute Eq. (3) by marching along each rendering ray in the reverse order $q = R, R-1, \dots$, since the vectors \mathbf{a}_j are accumulated from sample q onwards, and the opacity $f_\sigma(\mathbf{x}_q)$ affects only the visibility of successive samples $\mathbf{x}_q, \mathbf{x}_{q+1}, \dots$. To make this possible, we cache the final transmittance T_R , which is computed in the forward pass (this amounts to one scalar per ray). In backpropagation, we sequentially compute $\sigma(x_j)$ for every 3D point along the ray, and calculate $T_{j-1} = T_j \cdot \exp(\Delta\sigma(x_j))$ from T_j . This way, similar to the forward pass, the kernel only stores the accumulation of per-point features instead of keeping them all in memory.

Difference from checkpointing. Note that the latter is very different from “checkpointing” which can be trivially enabled for the naive renderer implementation in autograd frameworks such as PyTorch. This is because, unlike our memory-efficient backward pass from Eq. (3), a checkpointed backward pass still entails storing all intermediate features along rendering rays in memory.

4 Example applications

We show various 3D applications that could be boosted by the proposed components, from single-scene optimization with image-level losses to a versatile framework for large-scale 3D reconstruction and generation. Results are in Sec. 5.

Single-scene optimization with image-level losses. Constrained by intensive memory usage during rendering, existing volumetric methods are limited to optimizing pixel-level losses on a subset of rays, such as MSE, or using image-level losses on low-resolution images (64×64). In contrast, we show how *Renderer* allows seamless usage of image-level losses on high-resolution renders.

Multi-view reconstruction. Combining *Renderer* and *Splatter*, we introduce a versatile pipeline for 3D reconstruction and generation. Given a set of views (viewset) $\mathcal{V} = \{I_i\}_{i=1}^N$ and corresponding cameras $\{\pi_i\}_{i=1}^N$, we train a large-scale model Φ , which directly outputs the 3D representations $\theta = \Phi(\mathcal{V}, \pi)$ of the corresponding scene by learning 3D priors from large-scale data. Reconstruction starts by extracting a pixel-wise feature map $\mathbf{v} = \psi(I_i)$ from each image I_i and lifting them into the 3D representation $\hat{\theta}$ with *Splatter*. Model Φ takes $\hat{\theta}$ as input and outputs the final 3D representations $\theta = \Phi_\theta(\hat{\theta})$. Finally, *Renderer* outputs novel view images $\hat{I} = \mathcal{R}(\theta, \pi)$ from θ , and the model is trained by minimizing the loss \mathcal{L} between the novel rendered image and the corresponding ground truth I .

3D generation using viewset diffusion. Following recent works [2, 79], this 3D reconstruction pipeline could be extended into a diffusion-based 3D generator with very few changes. This is achieved by considering a noised viewset as



Fig. 5: Single-scene optimization with image-level losses. The memory efficiency of *Lightplane* allows rendering high resolution images in a differentiable way and backpropagating image-level losses. We show pre-optimized 3D scenes (in unseen views) and their stylizations with perceptual losses.

input to the network, and training the model to denoise the viewset, where each image I_i is replaced with $I_{it} = \alpha_t I_i + \sigma_t \epsilon_i$ where t is the noising schedule, α_t and $\sigma_t = \sqrt{1 - \alpha_t^2}$ are the noise level, and ϵ_i is a random Normal noise vector. During inference, the model initializes the viewset with Gaussian noise and iteratively denoises by applying the reconstruction model. This process simultaneously generates multiple views of the object as well as its 3D model.

5 Experiments

We first benchmark the performance of proposed components, and then demonstrate their practical usage for various 3D tasks, including single-scene optimization with image-level loss, and boosting the scalability of large-scale 3D models.

The scalability boost comes from both *input-size* and *modeling*. For input-size, it dramatically increases the amount of 2D information lifted to 3D by enlarging the number of input views and the output size. For modeling, the memory savings allow increasing the model and batch size during training.

5.1 Memory & speed benchmark

We measure components’ speed and memory in Figure 4. *Renderer* (left col.) is tested on a triplane with 256 points per ray, and compared to a PyTorch Autograd triplane renderer, adopted from [7, 9]. It easily supports high image sizes with low memory usage, which is unaffordable for the Autograd renderer. *Splatter* (right col.) is tested on lifting N input feature maps into a 160^3 voxel grid. We benchmark it against the lifting operations from PixelNeRF [98] and GeNVS [11], disabling the MLPs in *Splatter* for a fair comparison. As shown, *Splatter* can handle over a hundred views efficiently, while existing methods are restricted to just a few views. Speed-wise, both components are comparable to their autograd counterparts. See supplementary material for more results.

5.2 Single-scene Optimization with Image-level Loss.

The memory efficiency of the proposed components, in particular *Renderer*, allows rendering high-resolution images (*e.g.* 512^2) in a differentiable way with little memory overhead. Therefore, we can seamlessly use models which take full

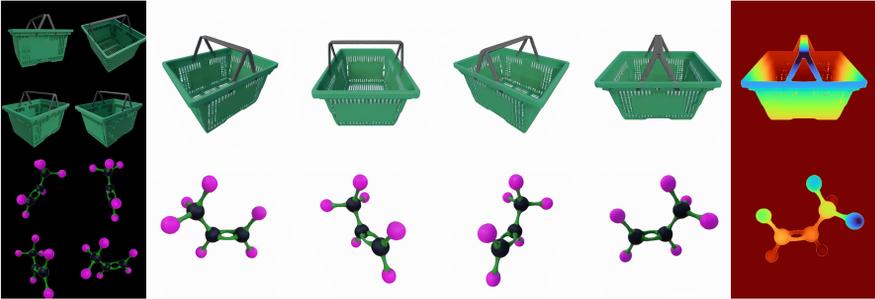


Fig. 6: Multi-view Large Reconstruction Model (LRM) with *Lightplane*. Taking four views as input (leftmost column), we show the RGB renders (mid) and depth (rightmost column) of the 3D reconstruction.

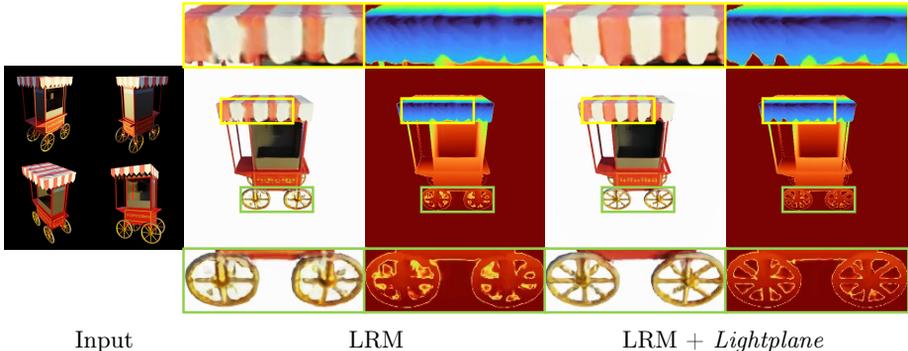


Fig. 7: Visual Comparison of LRM. Adding *Lightplane* to LRM gives more accurate geometry and appearance with little additional computation and memory cost.

images as input for loss calculation, *e.g.* perceptual loss [33], LPIPS [102], or SDS [61], and backpropagate these losses back to neural 3D fields. Constrained by memory usage, existing methods are limited to very low-resolution rendering [47, 61] or complicated and inefficient deferred backpropagation [100], while *Lightplane* can handle image-level losses easily. We take neural 3D field stylization as an example in Figure 5 and discuss more applications in supplementary.

Table 1: Quantitative results of LRM. The proposed method could effectively improve the reconstruction results, especially geometry (depth L1).

Method	PSNR \uparrow	LPIPS \downarrow	IOU \uparrow	Depth L1 \downarrow
LRM [31]	23.7	0.113	0.904	0.208
<i>Lightplane</i> +LRM	24.1	0.106	0.916	0.168



Fig. 8: Monocular 3D Reconstruction. With a single clean image as input (1st col.), our model could generate realistic 3D structures matching the input.

5.3 Multi-view LRM with *Lightplane*

We first validate *Lightplane*'s efficacy, in particular of the triplane *Renderer* and *Spatter*, by combining *Lightplane* with Large Reconstruction Model(LRM) [31]. Taking four images as input, this model outputs triplane as the 3D representation



Fig. 9: Unconditional 3D Generation displaying samples from our *Lightplane*-augmented Viewset Diffusion trained on CO3Dv2 [67].

Table 2: Amortized 3D Reconstruction. Our feedforward reconstructor (*Lightplane*) trained on the *whole* CO3Dv2 significantly outperforms baseline ViewFormer [41]. We further compare overfitting baselines (Voxel, NeRF [54]) to *Lightplane*, and to their scene-tuned versions (“Feedforward + Overfit”). Initializing from *Lightplane*-feedforward removes defective geometry leading to better depth error.

Method	Mode	#views	PSNR \uparrow	LPIPS \downarrow	Depth corr. \uparrow	Time \downarrow
ViewFormer [41]	Feedforward	9	16.4	0.274	N/A	N/A
<i>Lightplane</i>	Feedforward	10	20.7	0.141	0.356	1.6 sec
<i>Lightplane</i>	Feedforward	20	20.9	0.136	0.382	1.9 sec
<i>Lightplane</i>	Feedforward	40	21.4	<u>0.131</u>	<u>0.405</u>	2.5 sec
<i>Lightplane</i>	Feedforward + Overfit	160	<u>26.2</u>	0.086	0.449	5 min
Voxel	Overfit from scratch	160	26.5	0.086	0.373	35 min
NeRF	Overfit from scratch	160	26.3	0.108	0.658	1 day

via a series of transformer blocks. Every 3 transformer blocks (*i.e.* 5 blocks in total), we insert the *Splatter* layer, which splats source view features into a new triplane, taking previous block outputs as prior shape encoding. Plugging *Lightplane* into LRM adds little computational overheads, while clearly improving the performance. Additionally, the memory efficiency of our renderer enables LPIPS optimization without the added complexity of the deferred backpropagation in LRM [31]. We show the results in Table 1 and Figure 6, 7.

5.4 Large-Scale 3D Reconstruction and Generation.

Datasets and Baselines. We use CO3Dv2 [67] as our primary dataset, a collection of real-world videos capturing objects across 51 common categories. We implement the versatile model for 3D reconstruction and generation as described in Sec. 4, and extend *Lightplane* to unbounded scenes by contracting the ray-point’s coordinates [5] to represent background. Without loss of generality, we utilize UNet [69] with attention layers [84] to process 3D hash structures.

Amortized 3D Reconstruction. Existing amortized 3D reconstruction and novel view synthesis methods [67, 83, 88, 98, 103] only consider a few views (up to 10) as input due to memory constraints. Here, we enlarge the number of input views significantly. Unlike existing category-specific models, we train a *single*

Table 3: Unconditional 3D Generation on CO3Dv2. Our *Lightplane* significantly outperforms HoloDiffusion [37] and Viewset Diffusion [79]. It even beats HoloFusion [35], a distillation-based method, which takes 30 mins for one generation.

Method	Feed-forward	Hydrant		Teddybear		Apple		Donut		Mean		Inference Time
		FID ↓	KID ↓									
HoloFusion [35]	×	66.8	0.047	87.6	0.075	69.2	0.063	109.7	0.098	83.3	0.071	30mins
HoloDiffusion [37]	✓	100.5	0.079	109.2	0.106	94.5	0.095	115.4	0.085	122.5	0.102	<2min
Viewset Diffusion [79]	✓	150.5	0.124	219.7	0.178	-	-	-	-	-	-	<2min
<i>Lightplane</i>	✓	<u>75.1</u>	<u>0.058</u>	87.9	0.070	32.6	0.019	44.0	0.019	59.9	0.042	<2min

model on *all* CO3Dv2 categories, targeting a universal reconstruction model that can work on a variety of object types, and provide useful 3D priors for the following 3D optimizations. During training, 20 source images from a training scene are taken as inputs and MSE losses are calculated on five other novel views.

We evaluate in two regimes: (1) comparing our model to other feedforward baselines and single-scene overfitting methods to evaluate the model’s performance; (2) finetuning feedforward results using training views in a single scene to show the efficacy of our model as a learned 3D prior. Since few generalizable NeRF methods can work on all categories, we take ViewFormer [39] as the feedforward model baseline, which directly outputs novel view images using Transformer. In (2), we use 80 views as inputs to the feedforward model for initialization and report results of vanilla NeRF [54], and voxel-grid overfits (*i.e.*, trained from scratch). We evaluate results on novel views of unseen scenes.

Our model generates compelling reconstructions with just a single forward pass, shown in Tab. 2. After fine-tuning, it is on par with the overfitting baselines in color accuracy (PSNR, LPIPS), but largely outperforms the hash-based baselines (Voxel) in depth error. Since the frames of CO3D’s real test scenes exhibit limited viewpoint coverage, overfitting with hashed representations leads to strong defects in geometry (see Supp.). Here, by leveraging the memory-efficient *Lightplane* for pre-training on a large dataset, our model learns a generic surface prior which facilitates defect-free geometry. NeRF is superior in depth error while being on par in PSNR, at the cost of $\sim 50\times$ longer training time.

Unconditional Generation. Our model is capable of unconditional generation with only minor modifications, specifically accepting noisy input images and rendering the clean images through a denoising process. Utilizing the *Splatter* and *Renderer*, we can denoise multiple views (10 in experiments) within each denoising iteration, which significantly enhances the stability of the process and leads to markedly improved results. In the inference stage, we input 10 instances of pure noise and proceed with 50 Denoising Diffusion Implicit Model (DDIM) [75] sampling steps. We compare our method to Viewset Diffusion [79] and HoloFusion [35] quantitatively in Tab. 3 and evaluate qualitatively in Fig. 9. Our results significantly outperform other feedforward generation models and are comparable to distillation-based method, which is very time-consuming.

Conditional Generation. We can also introduce one clean image as conditioning, enabling single-view reconstruction. Moreover, our framework is also

amenable to extension as a text-conditioned model, utilizing captions as inputs. We show results and comparison in Figure 17 and Supp.

6 Conclusion

We have introduced *Lightplane*, a versatile framework that provide two novel components, *Splatter* and *Renderer*, which address the key memory bottleneck in network that manipulate neural fields. We have showcased the potential of these primitives in a number of applications, boosting models for reconstruction, generation and more. Once released to the community, we hope that these primitives will be used by many to boost their own research as well. ³

7 Acknowledgement

This work was done during Ang Cao’s internship at Meta AI as well as at the University of Michigan, partially supported by a grant from LG AI Research. We thank Roman Shapovalov, Jianyuan Wang, and Mohamed El Banani for their valuable help and discussions.

³ We discuss limitations and potential negative impact in Supp.

References

1. Anciukevicius, T., Xu, Z., Fisher, M., Henderson, P., Bilen, H., Mitra, N.J., Guerrero, P.: RenderDiffusion: Image diffusion for 3d reconstruction, inpainting and generation. arXiv.cs **abs/2211.09869** (2022) [4](#)
2. Anciukevicius, T., Xu, Z., Fisher, M., Henderson, P., Bilen, H., Mitra, N.J., Guerrero, P.: Renderdiffusion: Image diffusion for 3d reconstruction, inpainting and generation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12608–12618 (2023) [10](#)
3. Anonymous: Instant3d: Fast text-to-3d with sparse-view generation and large reconstruction model. Under Review (2023) [2](#), [23](#)
4. Barron, J.T., Mildenhall, B., Tancik, M., Hedman, P., Martin-Brualla, R., Srinivasan, P.P.: Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 5855–5864 (2021) [3](#)
5. Barron, J.T., Mildenhall, B., Verbin, D., Srinivasan, P.P., Hedman, P.: Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5470–5479 (2022) [3](#), [13](#), [22](#), [24](#)
6. Barron, J.T., Mildenhall, B., Verbin, D., Srinivasan, P.P., Hedman, P.: Zip-nerf: Anti-aliased grid-based neural radiance fields. arXiv preprint arXiv:2304.06706 (2023) [4](#)
7. Cao, A., Johnson, J.: Hexplane: A fast representation for dynamic scenes. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 130–141 (2023) [6](#), [11](#), [23](#)
8. Cao, Z., Hong, F., Wu, T., Pan, L., Liu, Z.: Large-vocabulary 3d diffusion model with transformer. arXiv preprint arXiv:2309.07920 (2023) [23](#)
9. Chan, E.R., Lin, C.Z., Chan, M.A., Nagano, K., Pan, B., De Mello, S., Gallo, O., Guibas, L.J., Tremblay, J., Khamis, S., et al.: Efficient geometry-aware 3d generative adversarial networks. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 16123–16133 (2022) [3](#), [4](#), [6](#), [11](#)
10. Chan, E.R., Lin, C.Z., Chan, M.A., Nagano, K., Pan, B., Mello, S.D., Gallo, O., Guibas, L., Tremblay, J., Khamis, S., Karras, T., Wetzstein, G.: Efficient geometry-aware 3D generative adversarial networks. In: arXiv (2021) [6](#)
11. Chan, E.R., Nagano, K., Chan, M.A., Bergman, A.W., Park, J.J., Levy, A., Aitala, M., Mello, S.D., Karras, T., Wetzstein, G.: GeNVs: Generative novel view synthesis with 3D-aware diffusion models. In: ICCV (2023) [2](#), [4](#), [5](#), [7](#), [9](#), [11](#), [22](#)
12. Chen, A., Xu, Z., Geiger, A., Yu, J., Su, H.: Tensorf: Tensorial radiance fields. In: European Conference on Computer Vision (ECCV) (2022) [2](#)
13. Chen, A., Xu, Z., Geiger, A., Yu, J., Su, H.: TensorRF: Tensorial radiance fields. In: arXiv (2022) [3](#), [26](#)
14. Chen, A., Xu, Z., Geiger, A., Yu, J., Su, H.: Tensorf: Tensorial radiance fields. In: European Conference on Computer Vision (ECCV) (2022) [6](#)
15. Chen, A., Xu, Z., Zhao, F., Zhang, X., Xiang, F., Yu, J., Su, H.: Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 14124–14133 (2021) [2](#)
16. Choy, C.B., Xu, D., Gwak, J., Chen, K., Savarese, S.: 3d-r2n2: A unified approach for single and multi-view 3d object reconstruction. In: Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part VIII 14. pp. 628–644. Springer (2016) [3](#)

17. Dao, T., Fu, D., Ermon, S., Rudra, A., Ré, C.: Flashattention: Fast and memory-efficient exact attention with io-awareness. *Advances in Neural Information Processing Systems* **35**, 16344–16359 (2022) [9](#)
18. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: 2009 IEEE conference on computer vision and pattern recognition. pp. 248–255. Ieee (2009) [24](#)
19. Dhariwal, P., Nichol, A.: Diffusion models beat gans on image synthesis. *Advances in neural information processing systems* **34**, 8780–8794 (2021) [4](#)
20. Fan, H., Su, H., Guibas, L.J.: A point set generation network for 3d object reconstruction from a single image. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 605–613 (2017) [3](#)
21. Fridovich-Keil, S., Meanti, G., Warburg, F.R., Recht, B., Kanazawa, A.: K-planes: Explicit radiance fields in space, time, and appearance. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12479–12488 (2023) [23](#)
22. Fridovich-Keil, S., Yu, A., Tancik, M., Chen, Q., Recht, B., Kanazawa, A.: Plenoxels: Radiance fields without neural networks. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5501–5510 (2022) [2](#)
23. Gadelha, M., Maji, S., Wang, R.: 3D shape induction from 2D views of multiple objects. In: arXiv (2016) [4](#)
24. Gao, J., Shen, T., Wang, Z., Chen, W., Yin, K., Li, D., Litany, O., Gojcic, Z., Fidler, S.: GET3D: A generative model of high quality 3d textured shapes learned from images. arXiv.cs **abs/2209.11163** (2022) [4](#)
25. Girdhar, R., Fouhey, D.F., Rodriguez, M., Gupta, A.: Learning a predictable and generative vector representation for objects. In: Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part VI 14. pp. 484–499. Springer (2016) [3](#)
26. Gkioxari, G., Johnson, J., Malik, J.: Mesh R-CNN. In: Proc. ICCV (2019) [3](#)
27. Goodfellow, I.J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A.C., Bengio, Y.: Generative adversarial nets. In: Proc. NeurIPS (2014) [4](#)
28. Gu, J., Liu, L., Wang, P., Theobalt, C.: Stylenerf: A style-based 3d-aware generator for high-resolution image synthesis. arXiv preprint arXiv:2110.08985 (2021) [4](#)
29. Henzler, P., Mitra, N.J., Ritschel, T.: Escaping plato’s cave using adversarial training: 3D shape from unstructured 2D image collections. In: Proc. ICCV (2019) [4](#)
30. Henzler, P., Reizenstein, J., Labatut, P., Shapovalov, R., Ritschel, T., Vedaldi, A., Novotny, D.: Unsupervised learning of 3d object categories from videos in the wild. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2021) [4](#)
31. Hong, Y., Zhang, K., Gu, J., Bi, S., Zhou, Y., Liu, D., Liu, F., Sunkavalli, K., Bui, T., Tan, H.: Lrm: Large reconstruction model for single image to 3d. arXiv preprint arXiv:2311.04400 (2023) [2](#), [3](#), [4](#), [12](#), [13](#)
32. Irshad, M.Z., Zakharov, S., Liu, K., Guizilini, V., Kollar, T., Gaidon, A., Kira, Z., Ambrus, R.: Neo 360: Neural fields for sparse view synthesis of outdoor scenes. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 9187–9198 (2023) [4](#)

33. Johnson, J., Alahi, A., Fei-Fei, L.: Perceptual losses for real-time style transfer and super-resolution. In: Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part II 14. pp. 694–711. Springer (2016) [12](#)
34. Kar, A., Häne, C., Malik, J.: Learning a multi-view stereo machine. *Advances in neural information processing systems* **30** (2017) [4](#)
35. Karnewar, A., Mitra, N.J., Vedaldi, A., Novotny, D.: Holofusion: Towards photo-realistic 3d generative modeling. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 22976–22985 (2023) [4](#), [7](#), [14](#)
36. Karnewar, A., Ritschel, T., Wang, O., Mitra, N.: Relu fields: The little non-linearity that could. In: ACM SIGGRAPH 2022 Conference Proceedings. pp. 1–9 (2022) [3](#)
37. Karnewar, A., Vedaldi, A., Novotny, D., Mitra, N.: Holodiffusion: Training a 3D diffusion model using 2D images. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (2023) [4](#), [7](#), [14](#)
38. Kerbl, B., Kopanas, G., Leimkühler, T., Drettakis, G.: 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics* **42**(4) (July 2023), <https://repo-sam.inria.fr/fungraph/3d-gaussian-splatting/> [2](#), [3](#)
39. Kulhánek, J., Derner, E., Sattler, T., Babuška, R.: ViewFormer: NeRF-free neural rendering from few images using transformers. In: Proc. ECCV (2022) [14](#)
40. Kulhanek, J., Sattler, T.: Tetra-nerf: Representing neural radiance fields using tetrahedra. arXiv preprint arXiv:2304.09987 (2023) [3](#)
41. Kulh’aneK, J., Derner, E., Sattler, T., Babuvska, R.: Viewformer: Nerf-free neural rendering from few images using transformers. In: European Conference on Computer Vision (2022), <https://api.semanticscholar.org/CorpusID:247593819> [13](#)
42. Li, J., Tan, H., Zhang, K., Xu, Z., Luan, F., Xu, Y., Hong, Y., Sunkavalli, K., Shakhnarovich, G., Bi, S.: Instant3d: Fast text-to-3d with sparse-view generation and large reconstruction model. arXiv preprint arXiv:2311.06214 (2023) [4](#)
43. Li, J., Li, D., Savarese, S., Hoi, S.: Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. arXiv preprint arXiv:2301.12597 (2023) [26](#), [28](#)
44. Li, R., Gao, H., Tancik, M., Kanazawa, A.: Nerfacc: Efficient sampling accelerates nerfs. arXiv preprint arXiv:2305.04966 (2023) [25](#)
45. Li, R., Tancik, M., Kanazawa, A.: NerfAcc: A general nerf acceleration toolbox. arXiv.cs [abs/2210.04847](#) (2022) [25](#)
46. Li, Z., Müller, T., Evans, A., Taylor, R.H., Unberath, M., Liu, M.Y., Lin, C.H.: Neuralangelo: High-fidelity neural surface reconstruction. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8456–8465 (2023) [3](#)
47. Lin, C.H., Gao, J., Tang, L., Takikawa, T., Zeng, X., Huang, X., Kreis, K., Fidler, S., Liu, M.Y., Lin, T.Y.: Magic3d: High-resolution text-to-3d content creation. arXiv preprint arXiv:2211.10440 (2022) [12](#)
48. Lin, C., Gao, J., Tang, L., Takikawa, T., Zeng, X., Huang, X., Kreis, K., Fidler, S., Liu, M., Lin, T.: Magic3D: High-resolution text-to-3d content creation. arXiv.cs [abs/2211.10440](#) (2022) [26](#)
49. Liu, R., Wu, R., Van Hoorick, B., Tokmakov, P., Zakharov, S., Vondrick, C.: Zero-1-to-3: Zero-shot one image to 3d object. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 9298–9309 (2023) [22](#)

50. Liu, S., Li, T., Chen, W., Li, H.: Soft rasterizer: A differentiable renderer for image-based 3d reasoning. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 7708–7717 (2019) [2](#)
51. Lombardi, S., Simon, T., Saragih, J., Schwartz, G., Lehrmann, A., Sheikh, Y.: Neural volumes: Learning dynamic renderable volumes from images. arXiv preprint arXiv:1906.07751 (2019) [3](#)
52. Luo, T., Rockwell, C., Lee, H., Johnson, J.: Scalable 3d captioning with pretrained models. arXiv preprint arXiv:2306.07279 (2023) [28](#)
53. Melas-Kyriazi, L., Rupperecht, C., Vedaldi, A.: Pc² projection-conditioned point cloud diffusion for single-image 3d reconstruction (2023). <https://doi.org/10.48550/ARXIV.2302.10668>, <https://arxiv.org/abs/2302.10668> [4](#)
54. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. In: ECCV (2020) [2](#), [3](#), [5](#), [6](#), [13](#), [14](#), [26](#), [30](#)
55. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: NeRF: Representing scenes as neural radiance fields for view synthesis. In: Proc. ECCV (2020) [26](#)
56. Müller, T., Evans, A., Schied, C., Keller, A.: Instant neural graphics primitives with a multiresolution hash encoding. ACM Transactions on Graphics (ToG) **41**(4), 1–15 (2022) [2](#), [6](#), [25](#)
57. Müller, T., Evans, A., Schied, C., Keller, A.: Instant neural graphics primitives with a multiresolution hash encoding. In: Proc. SIGGRAPH (2022) [3](#), [25](#)
58. Nguyen-Ha, P., Karnewar, A., Huynh, L., Rahtu, E., Heikkilä, J.: Rgbd-net: Predicting color and depth images for novel views synthesis. In: Proceedings of the International Conference on 3D Vision (2021) [4](#)
59. Nguyen-Phuoc, T., Li, C., Theis, L., Richardt, C., Yang, Y.: HoloGAN: Unsupervised learning of 3D representations from natural images. arXiv.cs [abs/1904.01326](#) (2019) [4](#)
60. Niemeyer, M., Geiger, A.: Giraffe: Representing scenes as compositional generative neural feature fields. In: Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR) (2021) [4](#)
61. Poole, B., Jain, A., Barron, J.T., Mildenhall, B.: Dreamfusion: Text-to-3d using 2d diffusion. arXiv preprint arXiv:2209.14988 (2022) [2](#), [5](#), [12](#)
62. Poole, B., Jain, A., Barron, J.T., Mildenhall, B.: Dreamfusion: Text-to-3d using 2d diffusion. arXiv.cs [abs/2209.14988](#) (2022) [26](#)
63. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: International conference on machine learning. pp. 8748–8763. PMLR (2021) [26](#)
64. Reiser, C., Peng, S., Liao, Y., Geiger, A.: KiloNeRF: Speeding up neural radiance fields with thousands of tiny MLPs. arXiv.cs [abs/2103.13744](#) (2021) [3](#)
65. Reiser, C., Szeliski, R., Verbin, D., Srinivasan, P., Mildenhall, B., Geiger, A., Barron, J., Hedman, P.: Merf: Memory-efficient radiance fields for real-time view synthesis in unbounded scenes. ACM Transactions on Graphics (TOG) **42**(4), 1–12 (2023) [3](#)
66. Reizenstein, J., Shapovalov, R., Henzler, P., Sbordone, L., Labatut, P., Novotny, D.: Common Objects in 3D: Large-scale learning and evaluation of real-life 3D category reconstruction. In: Proc. CVPR (2021) [4](#)
67. Reizenstein, J., Shapovalov, R., Henzler, P., Sbordone, L., Labatut, P., Novotny, D.: Common objects in 3d: Large-scale learning and evaluation of real-life 3d cat-

- egory reconstruction. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 10901–10911 (2021) **13, 31, 32**
68. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models (2021) **24**
69. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5–9, 2015, Proceedings, Part III 18. pp. 234–241. Springer (2015) **13**
70. Rosu, R.A., Behnke, S.: Permutosdf: Fast multi-view reconstruction with implicit surfaces using permutohedral lattices. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8466–8475 (2023) **3**
71. Sajjadi, M.S.M., Meyer, H., Pot, E., Bergmann, U., Greff, K., Radwan, N., Vora, S., Lucic, M., Duckworth, D., Dosovitskiy, A., Uszkoreit, J., Funkhouser, T.A., Tagliasacchi, A.: Scene representation transformer: Geometry-free novel view synthesis through set-latent scene representations. CoRR **abs/2111.13152** (2021) **4**
72. Schwarz, K., Liao, Y., Niemeyer, M., Geiger, A.: Graf: Generative radiance fields for 3d-aware image synthesis. Advances in Neural Information Processing Systems **33**, 20154–20166 (2020) **4**
73. Sharma, P., Tewari, A., Du, Y., Zakharov, S., Ambrus, R., Gaidon, A., Freeman, W.T., Durand, F., Tenenbaum, J.B., Sitzmann, V.: Seeing 3d objects in a single image via self-supervised static-dynamic disentanglement. arXiv.cs **abs/2207.11232** (2022) **4**
74. Skorokhodov, I., Tulyakov, S., Wang, Y., Wonka, P.: Epigraf: Rethinking training of 3d gans. arXiv preprint arXiv:2206.10535 (2022) **4**
75. Song, J., Meng, C., Ermon, S.: Denoising diffusion implicit models. arXiv preprint arXiv:2010.02502 (2020) **14**
76. Sun, C., Sun, M., Chen, H.T.: Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) pp. 5449–5459 (2021), <https://api.semanticscholar.org/CorpusID:244477646> **26**
77. Sun, C., Sun, M., Chen, H.T.: Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5459–5469 (2022) **3**
78. Sun, C., Sun, M., Chen, H.: Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In: CVPR (2022) **6**
79. Szymanowicz, S., Rupperecht, C., Vedaldi, A.: Viewset diffusion:(0-) image-conditioned 3d generative models from 2d data. arXiv preprint arXiv:2306.07881 (2023) **4, 7, 10, 14, 22**
80. Tewari, A., Yin, T., Cazenavette, G., Rezchikov, S., Tenenbaum, J.B., Durand, F., Freeman, W.T., Sitzmann, V.: Diffusion with forward models: Solving stochastic inverse problems without direct supervision. In: arXiv (2023) **4**
81. Tillet, P., Kung, H.T., Cox, D.: Triton: an intermediate language and compiler for tiled neural network computations. In: Proceedings of the 3rd ACM SIGPLAN International Workshop on Machine Learning and Programming Languages. pp. 10–19 (2019) **3**
82. Tochilkin, D., Pankratz, D., Liu, Z., Huang, Z., Letts, A., Li, Y., Liang, D., Laforte, C., Jampani, V., Cao, Y.P.: Triposr: Fast 3d object reconstruction from a single image. arXiv preprint arXiv:2403.02151 (2024) **4**

83. Trevithick, A., Yang, B.: Grf: Learning a general radiance field for 3d scene representation and rendering. In: arXiv:2010.04595 (2020) **4**, **13**
84. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. *Advances in neural information processing systems* **30** (2017) **13**
85. Verbin, D., Hedman, P., Mildenhall, B., Zickler, T., Barron, J.T., Srinivasan, P.P.: Ref-nerf: Structured view-dependent appearance for neural radiance fields. In: 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 5481–5490. IEEE (2022) **3**
86. Wang, N., Zhang, Y., Li, Z., Fu, Y., Liu, W., Jiang, Y.G.: Pixel2mesh: Generating 3d mesh models from single rgb images. In: *Proceedings of the European conference on computer vision (ECCV)*. pp. 52–67 (2018) **3**
87. Wang, P., Liu, L., Liu, Y., Theobalt, C., Komura, T., Wang, W.: NeuS: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. arXiv.cs **abs/2106.10689** (2021) **3**
88. Wang, Q., Wang, Z., Genova, K., Srinivasan, P.P., Zhou, H., Barron, J.T., Martin-Brualla, R., Snavely, N., Funkhouser, T.A.: Ibrnet: Learning multi-view image-based rendering. 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) pp. 4688–4697 (2021), <https://api.semanticscholar.org/CorpusID:232045969> **2**, **13**
89. Wang, Q., Wang, Z., Genova, K., Srinivasan, P.P., Zhou, H., Barron, J.T., Martin-Brualla, R., Snavely, N., Funkhouser, T.A.: Ibrnet: Learning multi-view image-based rendering. In: *Proc. CVPR* (2021) **4**
90. Wang, T., Zhang, B., Zhang, T., Gu, S., Bao, J., Baltrusaitis, T., Shen, J., Chen, D., Wen, F., Chen, Q., et al.: Rodin: A generative model for sculpting 3d digital avatars using diffusion. arXiv preprint arXiv:2212.06135 (2022) **23**
91. Wang, Y., Han, Q., Habermann, M., Daniilidis, K., Theobalt, C., Liu, L.: NeuS2: Fast learning of neural implicit surfaces for multi-view reconstruction. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 3295–3306 (2023) **3**
92. Wu, C.Y., Johnson, J., Malik, J., Feichtenhofer, C., Gkioxari, G.: Multiview compressive coding for 3d reconstruction. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 9065–9075 (2023) **4**
93. Wu, S., Rupprecht, C., Vedaldi, A.: Unsupervised learning of probably symmetric deformable 3d objects from images in the wild. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 1–10 (2020) **4**
94. Xu, Q., Xu, Z., Philip, J., Bi, S., Shu, Z., Sunkavalli, K., Neumann, U.: PointNeRF: Point-based neural radiance fields. arXiv.cs **abs/2201.08845** (2022) **3**
95. Yang, G., Huang, X., Hao, Z., Liu, M.Y., Belongie, S., Hariharan, B.: Pointflow: 3d point cloud generation with continuous normalizing flows. In: *Proceedings of the IEEE/CVF international conference on computer vision*. pp. 4541–4550 (2019) **3**
96. Yariv, L., Gu, J., Kasten, Y., Lipman, Y.: Volume rendering of neural implicit surfaces. *Advances in Neural Information Processing Systems* **34**, 4805–4815 (2021) **3**
97. Yu, A., Fridovich-Keil, S., Tancik, M., Chen, Q., Recht, B., Kanazawa, A.: Plenoxels: Radiance fields without neural networks. arXiv preprint arXiv:2112.05131 (2021) **3**, **26**
98. Yu, A., Ye, V., Tancik, M., Kanazawa, A.: pixelnerf: Neural radiance fields from one or few images. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 4578–4587 (2021) **2**, **4**, **5**, **9**, **11**, **13**

99. Yu, Z., Peng, S., Niemeyer, M., Sattler, T., Geiger, A.: Monosdf: Exploring monocular geometric cues for neural implicit surface reconstruction. *Advances in neural information processing systems* **35**, 25018–25032 (2022) [3](#)
100. Zhang, K., Kolkin, N., Bi, S., Luan, F., Xu, Z., Shechtman, E., Snavely, N.: Arf: Artistic radiance fields (2022) [12](#)
101. Zhang, K., Riegler, G., Snavely, N., Koltun, V.: Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492* (2020) [3](#)
102. Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable effectiveness of deep features as a perceptual metric. In: *CVPR* (2018) [2](#), [5](#), [12](#)
103. Zhou, Z., Tulsiani, S.: Sparsefusion: Distilling view-conditioned diffusion for 3d reconstruction. *arXiv preprint arXiv:2212.00792* (2022) [13](#)

A Supplementary Videos

Please watch our attached video for a brief summary of the paper and more results. We include more generated results and 360-degree rendering videos from our reconstructed and generated 3D structures to show their 3D consistency.

B Social Impact

Our main contribution is *Lightplane Splatter* and *Renderer*, a pair of 3D components which could be used to significantly scale the mapping between 2D images and neural 3D fields. Beyond their integral role in our versatile pipeline for 3D reconstruction and generation, single scene optimization, and LRM with *Lightplane*, these components can also function as highly scalable plug-ins for various 3D applications. We earnestly hope that they will be instrumental in advancing future research.

Based on *Lightplane Splatter* and *Renderer*, we have established a comprehensive framework for 3D reconstruction and generation. Similar to many other generative models [[11](#), [49](#), [79](#)], it is important to note that the results generated by this framework have the potential to be used in the creation of synthetic media.

C Limitations & Discussions

Our motivation of introducing contract coordinates [[5](#)] is to assist the model in differentiating between foreground and background elements, thereby enhancing the quality of foreground generation and reconstruction. Although contract coordinates could represent unbounded scenes, our main focus is still on foreground objects, and reconstructing or generating unbounded backgrounds is beyond the scope of this paper. Therefore, we only sample limited points in unbounded regions, which leads to floaters, blurriness and clear artifacts in the background, as can be observed in videos. Also, generating diverse and realistic backgrounds is a challenging task and we leave it as a promising future direction.

Lightplane introduces a versatile approach for scaling the mapping between 2D and 3D in neural 3D fields, designed to be compatible with arbitrary 3D hash representations with differentiable sampling functions. While our validation of this design has focused on voxel and triplane models, its adaptability should allow for easy generalization to other 3D hash representations, such as Hash Table [3] or HexPlane [7,21]. We pick voxel grids and triplanes as their structures are easy to be processed by the existing neural networks while designing neural networks to process some other 3D hash structures like hash tables is still an open question. Developing neural networks to support other 3D hash structures is a promising direction to explore while beyond the scope of this paper.

Lightplane significantly solves the memory bottlenecks in neural 3D fields, making rendering and splatting a large number of images possible in the current 3D pipelines. Although *Lightplane* has comparable speed to existing methods, rendering and splatting a large number of images is still time-consuming, which may limit its utilization in real applications. For example, doing a forward and backward pass on 512×512 rendered images takes around 5 seconds for each iteration. For *Renderer*, the spent time grows linearly to the ray numbers when ray numbers are huge. Reducing the required time for large ray numbers would be a promising direction.

Sadly, we observe a performance gap between different 3D hash representations (*i.e.*, voxel grids and triplanes) in the versatile 3D reconstruction and generation framework. Without loss of generalization, we use 3D UNet to process voxel grids and 2D UNet to process Triplane. Three planes (XY, YZ, ZX) are concatenated into a single wide feature map and fed to 2DUNet. The self-attention mechanism is then applied across all patches from the three planes, making this network an extension of our 3DUNet designed for voxel grids. However, we observed that this neural network configuration does not yield flawless results. In 3D reconstruction tasks, the images rendered at novel viewpoints exhibit slight misalignments with the ground-truth images. For generative tasks, while the network can produce realistic samples, it occasionally generates flawed outputs that significantly impact the Fidelity (FID) and Kernel Inception Distance (KID) scores. Developing a more efficacious neural network model for TriPlane processing [3,8,90], which could effectively communicate features from three planes, presents a promising avenue for future research.

D *Lightplane* Details

D.1 Implementation Details

Normalization Process in *Splatter*. Starting from a zero-initialized hash θ , *Splatter* is done by accumulating v_{ij} to the hash cell (*i.e.* voxel grids or triplanes) that contain x_{ij} , using the same trilinear/bilinear weights used in the *Renderer* operator to sample θ . After accumulating over all M rays, each hash cell is normalized by the sum of all splatting bi/trilinear weights landing in the cell. The normalization operation employed in our method, analogous to average pooling, averages the information splatted at identical positions in the

hash θ . This process guarantees that the magnitudes of the splatted features are comparable to those of the input view features, a factor that is beneficial for the learning process.

In the actual implementation, we execute the splatting process twice within the *Splatter* kernel. Initially, we splat the features of the input image into θ . Subsequently, a second set of weight maps is created, matching the spatial dimensions of the input image features, but with a feature of a single-scale: 1. These weight maps are then splatted into θ_{weight} . During the second splatting process within the *Splatter* kernel, we deactivate the Multilayer Perceptrons (MLPs) and suspend sampling from prior hash representations. This modification is implemented because our objective is to tally the frequency and weights of points being splatted into the same position within the hash representations, instead of learning to regress features. Finally, we get θ/θ_{weight} .

Performing the splatting operation twice inevitably results in additional time and memory overhead. In practice, θ_{weight} is relatively lightweight while θ is more memory-intensive. This is because they have the same spatial shape while θ_{weight} has a feature dimension of only 1. The normalization step θ/θ_{weight} , which is implemented in PyTorch, will cache the heavy θ , thereby increasing memory usage. We manually cache θ_{weight} to normalize gradients during backpropagation.

Experimental Details. We use $160 \times 160 \times 160$ voxel grids and 160×160 triplanes in our model. The input images are processed using a VAE-encoder [68] trained on the ImageNet dataset [18] and are converted into 32-dimensional feature vectors. Both the *Splatter* and *Renderer* components are equipped with 3-layer MLPs with a width of 64. Regarding training, we conduct 1000 iterations per epoch. The generative model is trained over 100 epochs, taking approximately 4 days, while the reconstruction model undergoes 150 epochs of training, lasting around 6 days, on a setup of 16 A100 GPUs, processing the entire Co3Dv2 dataset.

For *Splatter*, we sample 160 points along the ray. For *Renderer*, we sample 384 points along the ray, rendering 256×256 images. Instead of using original contract coordinates [5], we use a slightly different version which maps unbounded scenes into a $[-1, 1]$ cube.

$$CC(\mathbf{x}) = 0.5 * \begin{cases} a * \mathbf{x} & \|\mathbf{x}\| \leq 1 \\ \left((2 - a) * \left(1 - \frac{1}{\|\mathbf{x}\|} \right) + a \right) \left(\frac{\mathbf{x}}{\|\mathbf{x}\|} \right) & \|\mathbf{x}\| > 1 \end{cases} \quad (4)$$

We introduce a scale a to control the ratio between foreground and background regions, where the foreground regions are mapped to $[-a/2, a/2]$. As we are using explicit 3D hash, mapping foreground regions into larger regions would be helpful to represent details. When $a = 1$, it becomes the normal contract coordinates. We convert X, Y, Z axes into contract coordinates independently.

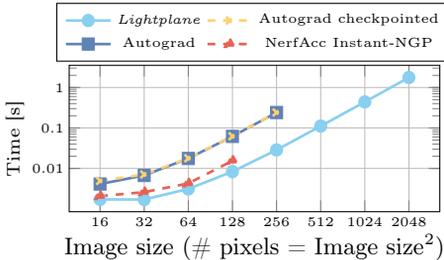


Fig. 1: Forward (FW) Time.

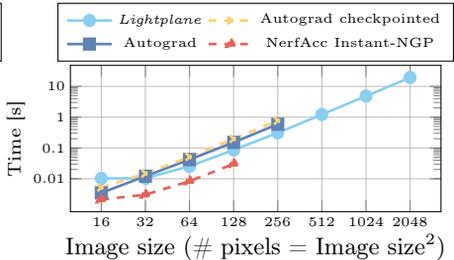


Fig. 2: Backward (BW) Time

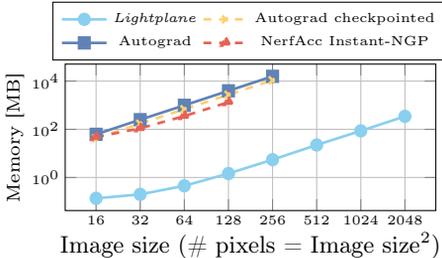


Fig. 3: Forward (FW) Memory.

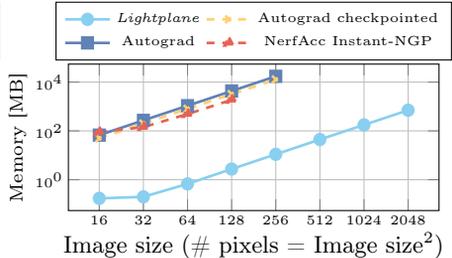


Fig. 4: Backward (BW) Memory

Fig. 10: *Lightplane Renderer* memory & speed benchmark showing the forward (FW and backward (BW) passes of *Lightplane Renderer*, compared to the *Autograd* renderer, *Checkpointing* (Pytorch checkpointing on *Autograd* renderer), and *NerfAcc Instant-NGP* (Instant-NGP [56] implemented in *NerfAcc* [44], which is claimed $1.1\times$ faster than the original version) *Lightplane* exhibits up to 4 orders of magnitude lower memory consumption at comparable speed. All axes are log-scaled.

D.2 *Lightplane* Performance Benchmark

Besides *Autograd* Renderer, implemented by pure Pytorch, we additionally compare *Lightplane Renderer* to two baselines: *Checkpointing* and *NerfAcc's Instant-NGP*, shown in Figure 10.

Checkpointing baseline applies the checkpointing technique in Pytorch to *Autograd* Renderer, which naive recalculates forward pass results during backward pass to save memories. Trivially applying checkpointing on *Autograd* indeed saves memories both in forward pass and backward pass, while still requires a large amount of memories, and cannot be used for large ray numbers.

NerfAcc's Instant-NGP is the Instant-NGP [57] implemented by *NerfAcc* [45], which is claimed to be $1.1\times$ faster than the original version of Instant-NGP, with tremendous optimization tricks for speed. Instant-NGP combines hash grid as 3D structures with fused MLP kernels (tiny-cuda-dnn), which is different from our *Renderer* with triplanes as 3D structures, and its internal settings are less flexible to change. To this end, it is hard to do a perfectly fair comparison. But still, we found that instant-NGP cannot work (will crash) with large image sizes,

Table 4: Quantitative results on NeRF Synthetic dataset [55].

Method	PSNR \uparrow	SSIM \uparrow	LPIPS $_{VGG}$ \downarrow
NeRF [54]	31.01	0.947	0.081
Plenoxels [97]	31.71	0.958	0.049
DVGO [76]	31.95	0.957	0.053
TensorRF-CP-384 [13]	31.56	0.949	0.076
TensorRF-VM-48 [13]	32.39	0.957	0.057
<i>Lightplane</i>	<u>32.12</u>	0.957	<u>0.050</u>

as they heavily rely on the L2 cache of GPUs for optimal speed, which memory is very limited and cannot support large image sizes. While their backward pass speed is significantly faster than *Lightplane Renderer*, it still cannot be extended to large output image sizes.

E More Results

E.1 Single Scene Optimization

Synthetic NeRF Results. We validate the correctness of *Lightplane* by overfitting on the Synthetic NeRF dataset, shown in Table 4. As the target is to show the convergence of *Lightplane*, we don’t employ any complicated tricks to optimize the performance and speed. *Lightplane* could get promising single-scene optimization results, demonstrating that it could be used as a reliable package in various 3D tasks.

DreamFusion with SDS Loss. The memory efficiency of *Lightplane* allows directly applying SDS [62] on high-resolution rendered images. As analyzed in Magic3D [48], existing 3D generations using SDS loss are limited to low-resolution rendered images: they first render low-resolution images for SDS to generate coarse 3D structures, and then convert the generated 3D structures into 3D meshes, which are used to generate high-resolution images. Using *Lightplane* allows direct optimization on high-resolution images.

Adversarial Attacking on LVM (Large Vision Model). We showcase another interesting application empowered by our *Lightplane* by adversarial attacking LVM models, *e.g.* CLIP [63] and BLIP2 [43] After rendering full images from the neural 3D field overfitted on a specific scene, we feed rendered images into CLIP model and calculate cosine similarity between image feature vectors and target text vectors, which similarity works as a loss to optimize the neural 3D fields.

E.2 Multi-view LRM with *Lightplane*

We show more results of Multi-view LRM with *Lightplane* in Figure 12 and Figure 13.



Fig. 11: 3D Adversarial Attacking on CLIP model. Given a fitted 3D scene (1st and 3rd column), we optimize the neural 3D fields so that features of rendered images are aligned to a specific text description, *i.e.* giraffe, in CLIP’s feature space, while keeping the appearance perceptually the same.

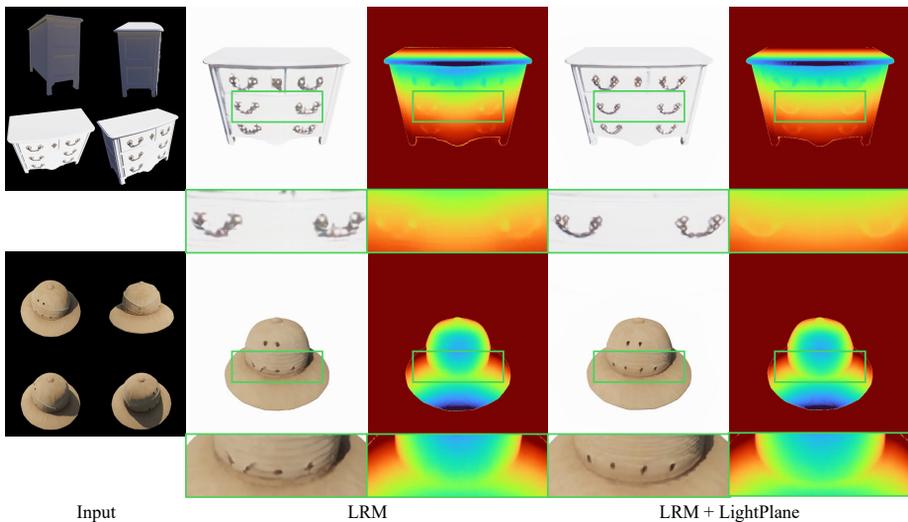


Fig. 12: Reconstruction comparison between LRM and LRM + *Lightplane*.

E.3 3D Reconstruction

We show amortized 3D reconstruction results after fine-tuning on a single scene in Figure 14, with voxel grids (*Lightplane-Vox*) and triplanes (*Lightplane-Tri*) as 3D structures. We compare them to overfitting results (training from scratch) using the 3D structures. Overfitting a single scene on Co3Dv2 dataset leads to defective 3D structures, like holes in depths. Initializing from the outputs of our amortized 3D reconstruction model could effectively solve this problem, leading to better results.

E.4 Unconditional Generation

We show 360-degree rendering for unconditional generation in Figure 15 and Figure 16.

E.5 Conditioned Generation

We show monocular 3D reconstruction with a single image as input in Figure 17, and text-conditioned generation in Figure 18. For text-conditioning experiments, we follow CAP3D [52]: we use BLIP2 [43] to generate captions of each image insides scenes and utilize LLAMA2 to output the comprehensive caption for the whole scene.



Fig. 13: Multi-view LRM with *Lightplane*.

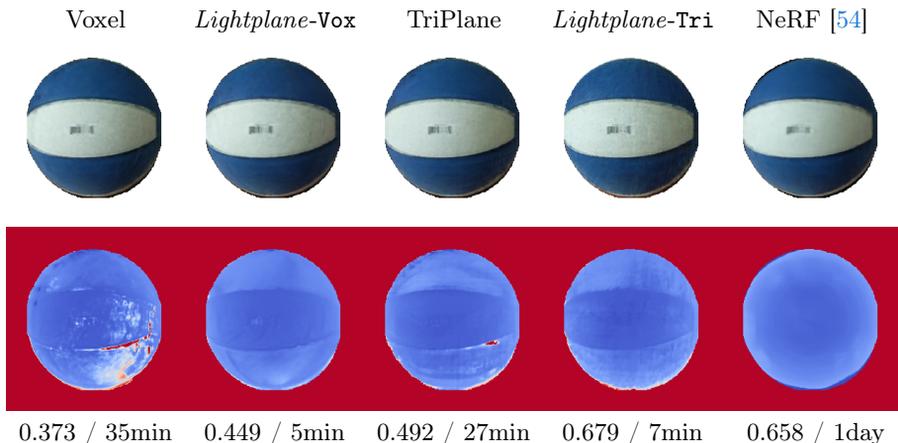


Fig. 14: 3D Reconstruction with Learned Initialization. We show rendered images (top row) and depths. Optimizing hashed representations (Voxel, Triplane) on real scenes leads to geometric defects. Using our models (*Lightplane-Vox*, *Lightplane-Tri*), we first learn a reconstruction prior on CO3Dv2. We then initialize reconstruction with a feed-forward pass accepting up to 100 source views of a single-scene. After fine-tuning, we observe improved quality of the reconstructed geometry (columns 3 and 4). We show *Depth Corr.* (\uparrow) and *Overfitting Time* (\downarrow) below images.

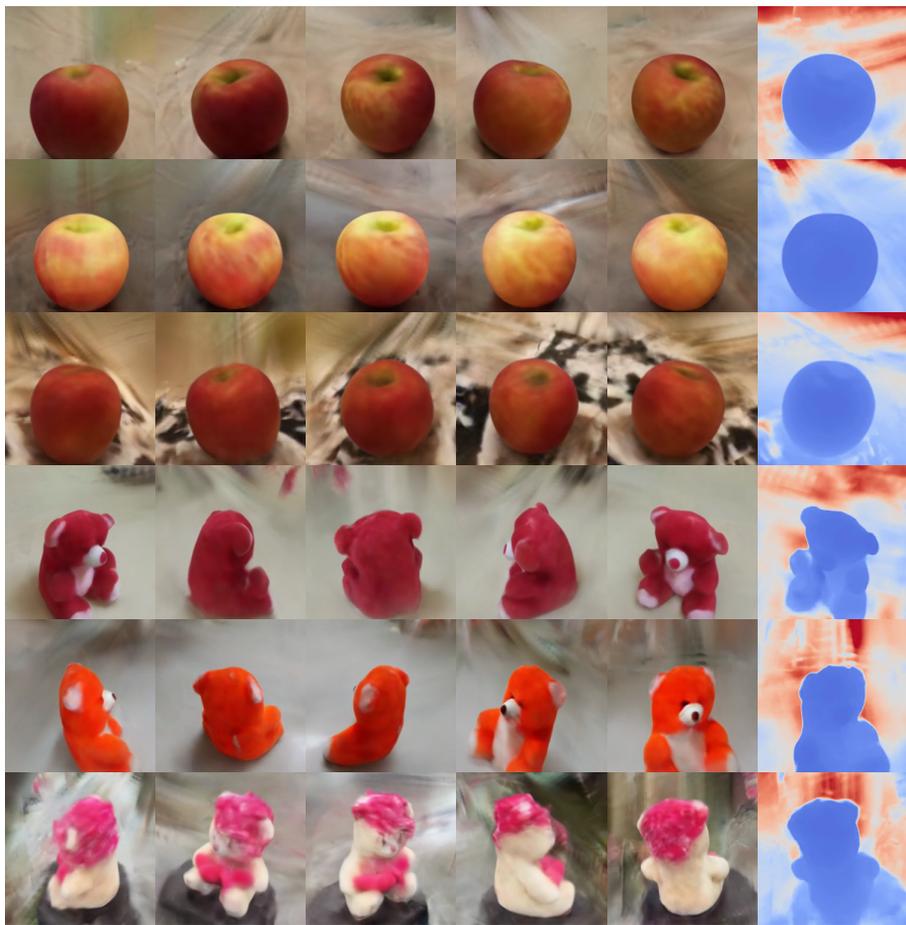


Fig. 15: Unconditional 3D Generation displaying samples from our *Lightplane*-augmented Viewset Diffusion trained on CO3Dv2 [67].

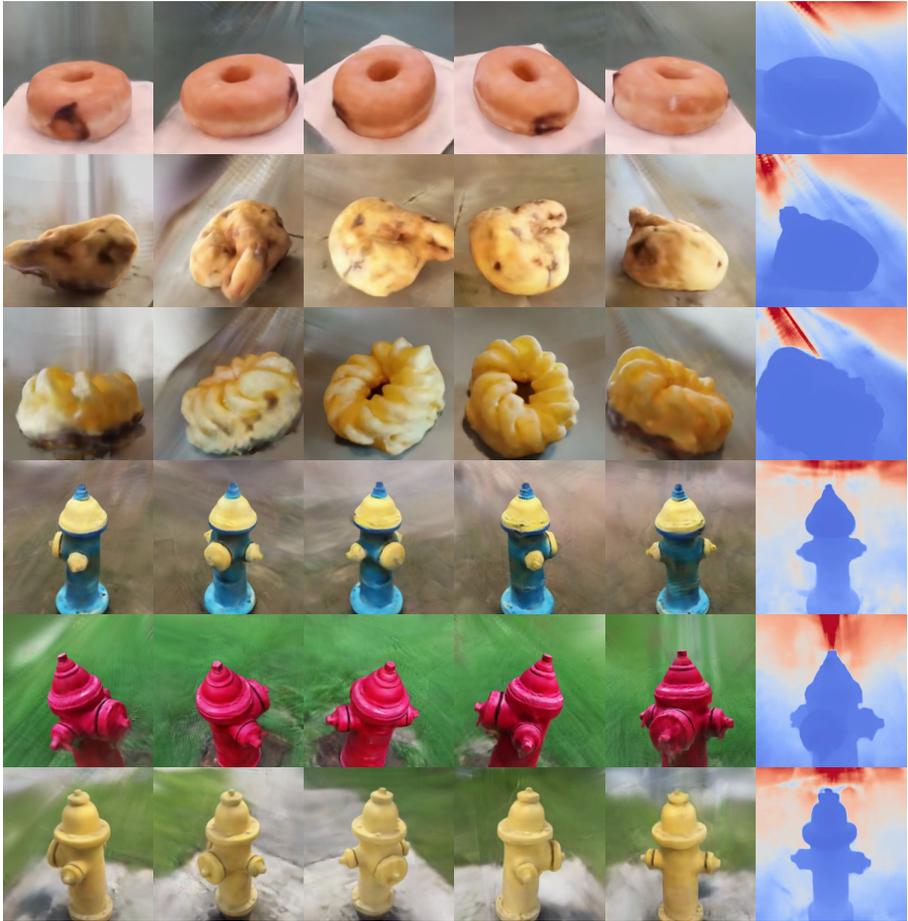


Fig. 16: Unconditional 3D Generation displaying samples from our *Lightplane*-augmented Viewset Diffusion trained on CO3Dv2 [67].



Fig. 17: Monocular 3D Reconstruction on CO3Dv2. With a single clean image as input, our model could generate realistic 3D structures matching the input views.



A white bowl with a blue and white fish in the center



A blue and white fire hydrant fire hydrant in a grassy area



A blue and white fire hydrant with a blue cap on the top

Fig. 18: Text-Conditioned Generation on CO3Dv2. Our pipeline could generate 3D structures with text input as conditions.