

Data Parallel Ray Tracing of Massive Scenes based on Neural Proxy

Shunkang Xu^{†1}, Xiang Xu^{‡2} , Yanning Xu¹  and Lu Wang^{‡1} 

¹School of Software, Shandong University, China.

²Shandong Key Laboratory of Blockchain Finance, Shandong University of Finance and Economics, China.

Abstract

Data-parallel ray tracing is an important method for rendering massive scenes that exceed local memory. Nevertheless, its efficacy is markedly contingent upon bandwidth owing to the substantial ray data transfer during the rendering process. In this paper, we advance the utilization of neural representation geometries in data-parallel rendering to reduce ray forwarding and intersection overheads. To this end, we introduce a lightweight geometric neural representation, denoted as a "neural proxy." Utilizing our neural proxies, we propose an efficient data-parallel ray tracing framework that significantly minimizes ray transmission and intersection overheads. Compared to state-of-the-art approaches, our method achieved a 2.29~3.36× speedup with an almost imperceptible image quality loss.

CCS Concepts

• **Computing methodologies** → **Computer graphics; Ray tracing; Neural networks;**

1. Introduction

With increasing demand for render quality of digital assets, scene data sizes increase significantly, even beyond the GPU or host memory. An important rendering method for such scenes is data-parallel (or data-distributed) rendering, whereby scene components are distributed across multiple compute units (e.g., HPC compute nodes or GPUs). Data-parallel ray tracing can provide high-quality rendering results for massive scenes. Nevertheless, such methods require rays to be sent to the node owning the relevant data (ray forwarding), which incurs expensive data transfer overheads between nodes.

The instancing technique [Car05] is commonly used to reduce the memory footprint of massive scenes when modeling and rendering, specifically for complex geometries that need to be duplicated multiple times. This technique utilizes transformation matrices to recycle the actual geometries, avoiding the repeated storage of the same mesh. Taking the coconut tree object in Disney's Moana Island scene [Stu18] as an example, a large coconut tree forest can be constructed using the base mesh (i.e., a single coconut tree mesh) and the affine transformation matrix of the instances.

Wald et al. [WP22] employed data-parallel ray tracing to render massive instanced scenes, where the scene is decomposed in object space and distributed across multiple nodes. In this framework, a top-level bounding volume hierarchy (BVH) is used to decide which node the ray should be sent to. However, the top-level BVH

is limited in providing precise geometric distribution information, resulting in scenarios where a ray needs to visit all nodes to find the closest hit. This results in the ray forwarding overhead often becoming a performance bottleneck, especially as the ray tracing capabilities of the GPU continue to be upgraded.

To reduce ray transmission overhead in data-parallel ray tracing, we aim to obtain precise intersection information of rays before the actual intersection. Consequently, we utilize two lightweight neural networks to represent the geometric distribution of each base mesh, termed a "neural proxy." In this way, any complex geometry can be compressed into a 1.19MB neural proxy. Considering the limited number of base meshes in a scene, we only need a small amount of additional memory to cache all neural proxies within each node. With our neural proxy, each secondary ray can be forwarded directly to the target node with the *closest-hits*, and the shadow ray can be computed locally without transmission. The main contributions include:

- a neural network-based geometry representation that can evaluate the geometry's intersection information at a specific ray origin and direction. This representation method is view and location-independent, once trained, can be used for any scene and different instances of the same scene, and supports rigid body transformations and model overlapping.
- a neural proxy-based data-parallel ray traversal method that can rapidly determine the intersection of the ray with individual instances in the scene without having to visit the actual geometries.
- an efficient data-parallel ray tracing framework for massive instanced scenes, where each node can determine the target node of each secondary ray locally, avoiding ineffective ray transmis-

[†] Shunkang Xu and Xiang Xu contributed equally to this work

[‡] Co-corresponding authors: Lu Wang (luwang_hcivr@sdu.edu.cn)

sion. At the same time, each node can locally generate ray-traced shadows for the entire scene without ray transmission.

2. Related work

2.1. Distributed Ray Tracing

Distributed ray tracing can be classified into image-parallel, data-parallel, and hybrid methods.

Image-parallel ray tracing involves partitioning the image space into numerous sub-regions assigned to individual compute nodes, with each node responsible for executing a portion of the ray task. The scene is dynamically loaded as required, avoiding data transmission among nodes. When ray tracing massive scenes that do not fit into local memory, this approach results in on-demand scene data fetching [DGBP05, KWR*17, PSL*98, WBS03]. Image-parallel ray tracing is frequently paired with Distributed Shared Memory (DSM) to hide the intricacy of memory scheduling. Jaros et al. [JRSŠ21] utilized the efficient and low-latency interconnections between NVIDIA GPUs to implement device memory sharing for massive scene rendering.

In data-parallel ray tracing, the scene is partitioned into sub-regions and then assigned to compute nodes. In this method, ray tasks need to be transferred among compute nodes. Kato et al. [KNE*01] evenly distributed the scene's primitives to different nodes. Each node separately intersects all rays and gathers the results in a master node to get the closest intersection. Fouladi et al. [FSP*22] applied data-parallel ray tracing on thousands of small, serverless cloud computing nodes. Their approach can render terabytes of scenes by appropriate scene data decomposition and resource scheduling. For single-node/multi-GPU hardware

with high-speed interconnections, Wald et al. [WJZ23] proposed a data-parallel ray-tracing framework that improves each GPU's computational and bandwidth utilization by allowing for redundant data transfers. For massive instanced scenes, Wald et al. [WP22] proposed an efficient data-parallel ray tracing architecture that reduces the ray-forwarding overhead by caching the instances' top-level BVHs at each node.

The hybrid approach aims to merge the image-parallel approach and the data-parallel approach, combining the respective advantages of each to improve the rendering efficiency. Reinhard et al. [RCJ99] introduced a multi-processor ray tracing algorithm with dynamic scheduling, wherein the relevant data is locally loaded for coherent rays proceeding, while the incoherent rays are sent to other processors owing the required data. Navrátil et al. [NFLC12, NCFL13] proposed various dynamic scheduling approaches for large-scale distributed nodes. Son et al. [SY17] proposed a timeline scheduling model for heterogeneous clusters of GPUs/CPU to increase resource utilization. Xu et al. [XWPG*24] employed the ray transmission information recorded in the previous frame to schedule the scene data and rays of the current frame during rendering, thus avoiding node synchronization during current frame rendering. Zellmann et al. [ZWB*22] proposed an *island parallelism* method that splits a set of $N \times M$ compute nodes into M islands of N nodes each and using data-parallel rendering within each island and image parallelism across islands.

2.2. Instancing Technique

The instancing technique is commonly used in 3D modeling and rendering to reduce storage overhead and GPU bandwidth consumption when transferring geometry [Car05]. Most ray trac-

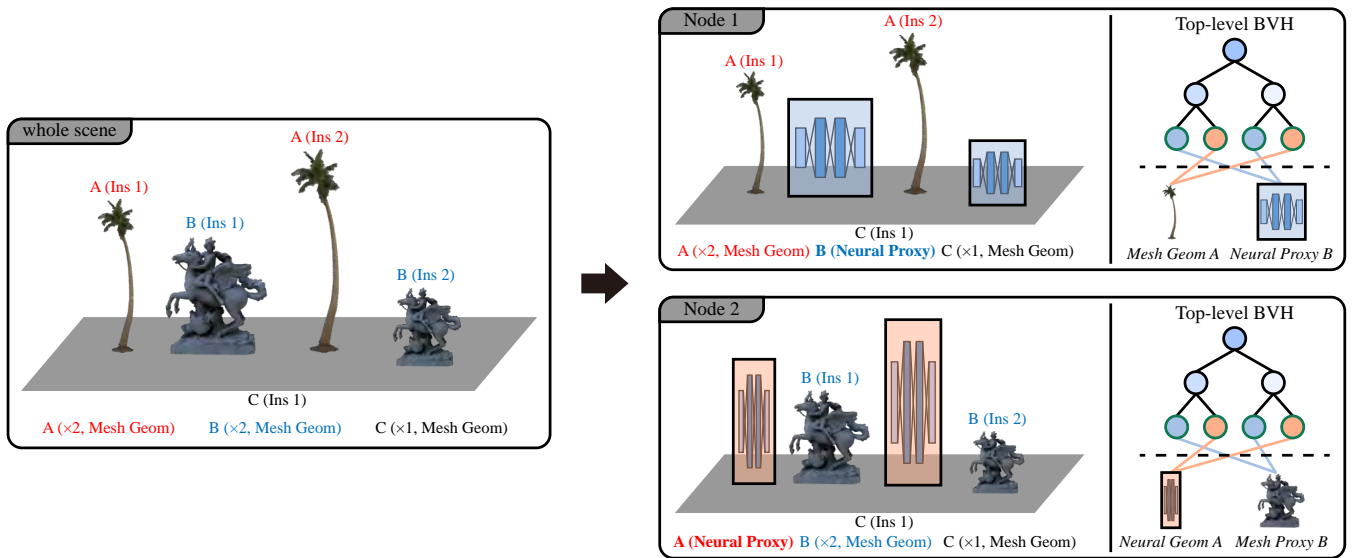


Figure 1: Illustration of data distribution of an instanced scene. On the left is an instanced scene constructed by three base meshes. There are two instances of mesh geometry A, two instances of mesh geometry B, and only one instance of mesh geometry C. The right shows the data distribution of two nodes; each node stores the top-level BVH, assigned local geometries, and neural proxies of other non-local geometries.

ing APIs and rendering systems support instancing [WWB*14, PBD*10, GIF*18]. Typical applications involve large outdoor environments that contain tree leaves and foliage. An early example of real-time ray tracing of a massively instanced model is that of Wald [Wal05]. Johansson et al. [Joh13] utilized the instancing technique to achieve real-time rendering of architectural models. Ivson et al. [IC14] showed the effectiveness of instancing when applied to large CAD models. By using a shape-matching algorithm, they reduced memory usage and achieved high rendering performance for massive 3D CAD models. Slomp et al. [SKD*13] used hardware-accelerated geometry instancing to render surfels and voxels efficiently.

2.3. Neural Implicit Geometry

Neural implicit geometry is a novel method for representing geometry. Many recent works [SMB*20] [LGZL*20] [MBRS*21] have utilized neural networks to represent scenes, with NeRF [MST*21] being the first to use the neural radiance field represented by a Multi-Layer Perceptron (MLP) for novel view synthesis. IDR [YKM*20] reconstructs surfaces using neural networks by representing geometry as the zero-level set of an MLP that is considered to be a signed distance field (SDF). Takikawa et al. [TLY*21] used an octree-based feature volume structure to adaptively fit shapes with multiple discrete LODs, enabling real-time rendering of high-fidelity neural SDFs.

Fujieda et al. [FKH23] first integrated the neural network-based approach into a BVH-based ray tracing pipeline. They proposed a neural intersection function that utilized two neural networks to fit the visibility information of the exterior and interior of an axis-aligned bounding box and then applied them to ray intersections to accelerate the rendering process instead of the actual intersection based on BVH. Subsequently, Weier et al. [WRM*24] employed hash grids [MESK22] for efficient 3D sparse data encoding, enabling representation of geometric information such as depth, normal, and appearance. These existing methods have successfully integrated neural implicit geometry into the ray tracing pipeline, resulting in high-quality rendering. However, the test scenes utilized in these methods are typically simple, and the neural implicit geometry has not yet been applied to large-scale scenes.

In this paper, we present, for the first time, a data-parallel ray tracing framework based on neural implicit geometry to render scenes that can not hold on a single node. In contrast to previous methods, we do not use neural geometry to directly replace the intersection results of real geometry. Instead, we employ our neural proxies to efficiently determine the distributed node to which a ray should be directed, thereby minimizing data transmission overhead. This approach allows us to preserve more geometric detail in large-scale scenes without the requirement of a complex and precise network. Meanwhile, given the substantial number of instances in "real" scenes, our neural proxy is built for the base geometry. Therefore, different instances of the same geometry can share a single neural proxy, reducing the training and inference overhead. Furthermore, the trained model can be directly applied to the same geometry in other scenes without re-training.

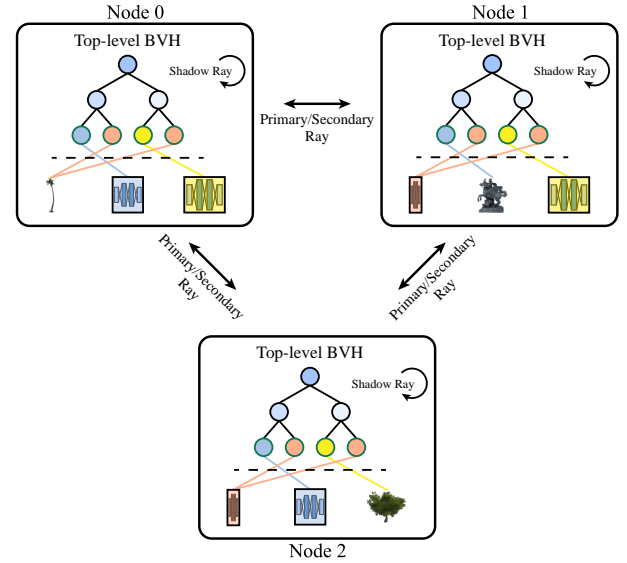


Figure 2: Our data-parallel ray tracing framework. Primary and secondary rays need to be transmitted among nodes, while shadow rays can be processed locally in each node.

3. Overview

In data-parallel ray tracing, the scene is partitioned into multiple parts and distributed across several nodes. Consequently, rays may need to be forwarded to all nodes to get the global *closest-hits*. The ray communication overhead often becomes a performance bottleneck for data-parallel ray tracing. Our data-parallel ray tracing framework employs neural proxies to predict which geometry will have the global *closest-hits* prior to the actual ray-geometry intersection. In this way, rays can be directly sent to the nodes containing the closest intersected geometry, thereby significantly reducing the overall ray transmission overhead.

3.1. Data Organization and Storage for Instanced Scenes

Geometric objects may be constructed by multiple levels of instancing. We choose a proper instancing level as the base mesh and construct the bounding box and its neural proxy. The standard for selecting a base mesh is that it contains a large number of meshes and is referenced multiple times in the scene. For example, the forest in Disney's Moana Island scene includes three instance levels (forest-tree-leaf), and we select the tree level as the base mesh in our practice.

We distribute the scene data in object space, where each node stores the assigned objects (each object contains the base mesh and its instances) as local data. Furthermore, each node needs to store the global information of the scene, including the top-level BVH (created over the bounding boxes of all instances in the scene) and neural proxies of non-local objects to determine the intersection of the rays with the entire scene, as shown in Fig. 1.

Our method requires each node to cache neural proxies for each non-local base mesh. Owing to the lightweight design of our neural

proxies, each one needs only approximately 1216 KB (1.19 MB) of storage. Consequently, even if there are 100 base meshes in the scene, each node only needs a few tens of MB of additional storage, which is almost negligible compared to the storage of the actual meshes.

3.2. Distributed Ray Tracing Framework

In our distributed ray tracing framework, we employed distinct methods to process primary, secondary, and shadow rays, illustrated in Fig.2. For the primary ray traversal, our approach is grounded in the work of Wald et al. [WP22]. Instead, secondary and shadow rays are traversed using our neural proxy-based traversal method, which can significantly reduce the transmission overhead of the secondary rays and allow shadow ray traversal to be performed directly and locally at the node. The detailed algorithm of our neural proxy-based ray traversal approach will be described in Section 5.

3.2.1. Primary Ray Traversal

We applied the wavefront ray traversal for primary rays. First, these rays get their global *closest-hits* across nodes by forwarding rays across nodes. Then, they re-traverse at the node containing their closest intersected geometry to complete shading. This strategy can significantly reduce the ray data that needs to be transmitted between nodes. Each ray to be transmitted only contains the information essential for traversal, including the origin, direction, throughput, pixel index, the current $tMax$ (the distance of the closest intersection), *next* (the index of the next instance that needs to be visited), etc.

During data-parallel traversal, we use the top-level BVH to minimize the number of ray transfers further. Firstly, the ray intersects with the BVH of local objects, updating the *next* and corresponding $tMax$. Secondly, the ray intersects with the top-level BVH to determine its intersection with non-local geometry. Finally, if the ray intersects with the non-local instance, it should be sent to the corresponding node for continued traversal, i.e., perform the first

two steps again. Otherwise, the ray will be sent to the corresponding node where the *next* object is located.

3.2.2. Secondary Ray Traversal

After the primary rays find the global *closest-hits* through distributed traversal, secondary rays are generated. It is notable that when we use the top-level BVH to determine which object the ray needs to be accessed, in the case of overlapping bounding boxes, a separate ray-geometry intersection is required for each overlapping object. This will result in a ray needing to be transmitted between nodes multiple times, particularly evident when the number of distributed nodes increases.

In our framework, each node stores neural proxies for non-local geometries to evaluate a precise intersection of a ray with all objects, avoiding redundant transmissions. Specifically, the secondary rays initially intersect with the BVH of local objects and then intersect with the top-level BVH. When a ray intersects a bounding box of a non-local object, we record the ray's hit information with the object's bounding box. The recorded information can subsequently be utilized to evaluate the intersection distance t of each non-local geometry using our neural proxies. This allows each node to predict locally which geometry has the global *closest-hits*. Then, the rays can be sent directly to the node that owns the target geometry to intersect and shade with the actual geometry.

3.2.3. Shadow Ray Traversal

For the shadow rays of our framework, given that the nodes can directly predict the ray intersection with all objects by neural proxies, the shadow ray traversal can proceed locally at each node without any ray transmission. Our shadow traversal method is similar to the secondary Ray Traversal, as it follows the sequence of intersecting rays with local objects, top-level BVHs, and neural proxies. The difference is that the shadow ray only needs to determine whether it is occluded by local objects or neural proxies. If the shadow ray is not occluded, the corresponding contribution is added to the frame buffer.

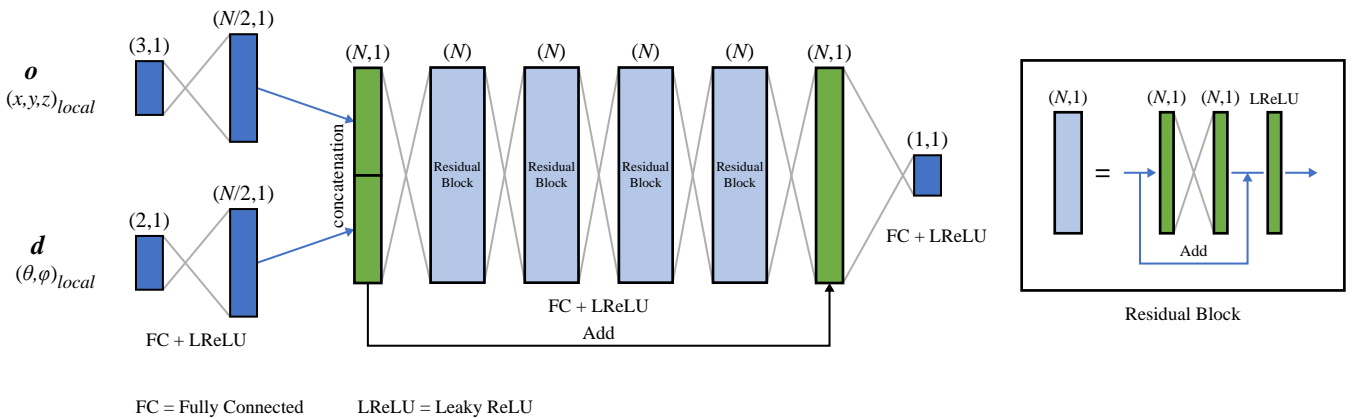


Figure 3: The network structure of our neural proxy with the residual blocks shown on the right. Our visibility and depth networks use the same structure and output visibility and t , respectively.

In summary, our data-parallel ray tracing framework can send the secondary rays directly to the node containing its *closest-hits*, regardless of the number of distributed nodes. Additionally, the shadow can be generated locally at each node without ray transmission. Also, our lightweight neural geometry design makes our occlusion evaluation faster than ray intersection with actual meshes.

4. Neural Proxy of Instanced Scene

We devised a lightweight neural geometric representation called neural proxy that described the geometric distribution in its bounding box. With our neural proxy, we can rapidly determine the intersection distance between a ray and a geometry without visiting its actual geometry. At the same time, the model supports affine transformations, which enable seamless integration with instancing techniques and improve the efficiency of our data-parallel framework in rendering massive instancing scenes.

Our choice to represent the geometric distribution with a neural network instead of other information such as radiance, material, etc., is based on two factors. First, the geometric distribution is independent of view and location, so it can be applied to multiple instances. Second, the inability to obtain the global closest intersection at a single node is a primary contributor to the communication overhead in data-parallel ray tracing, so we aim to determine the intersection of the ray with all scene geometries directly at each node.

Our neural proxy is constructed for the base mesh so that all instances of the same mesh can share the same neural network. Each neural proxy contains two distinct neural networks: visibility and depth networks. The visibility network is used for these rays with origins outside the bounding box, while the depth network is used together when a ray's origin is inside the bounding box. Both networks take the same input: position o located on the surface of the bounding box and direction d in local space. We use both networks to handle the various types of ray intersections in the scene, which are represented as $visibility(o, d, visibility)$ and $depth(o, d, t)$.

While a singular network may offer precise results for smooth meshes by simultaneously outputting both visibility and depth, our method is oriented to massive scenes, and most meshes have high-frequency contour information. Therefore, we consistently employ two distinct networks to learn visibility and depth information independently to enhance accuracy.

4.1. Network Structure

The visibility and depth network structures use the same network structure, as shown in Fig. 3. We utilize a fully connected neural network with residual blocks [ZBX*21]. The size N of layers is a hyperparameter that can be adjusted according to the geometric complexity. In our experiments, we set the size N for both the visibility and depth networks to 256 to balance the inference speed and rendering quality.

4.2. Data Preparation and Training

To generate the training data uniformly, we divided each face of the bounding box into 128×128 grids under the local space. The

position o is subsequently generated through random sampling on each grid. Concurrently, the hemispherical space is divided into eight regions, and random sampling is performed within each region to obtain the direction d , as shown in Fig. 4. Each pair of rays (o, d) intersected with the inner geometry using the OptiX pipeline [PBD*10] to acquire corresponding *visibility* and t values. Before training, we apply the Min-Max normalization to positions $o(x, y, z)$. The direction $d(\theta, \phi)$ on the hemisphere is also normalized to $[0, 1]$ by dividing them by π and 2π respectively. Visibility values remain unprocessed, while depth values require normalization by dividing them with the bounding box diagonal length.

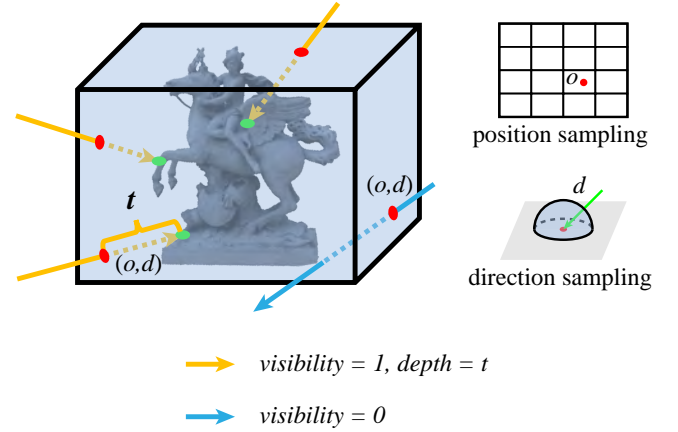


Figure 4: Training data generation. The rays are generated by sampling origin position o on the bounding box and direction d under local space. The visibility and depth t of rays are then obtained by intersection.

When training the visibility network, a specific ratio of ray data with visibility values of 0 and 1 is preserved in the dataset, and the MSE loss equation is employed. For the depth network, we used the standard L1 loss equation for training, and only ray data with valid depth values are preserved in the dataset to make the network concentrate on fitting the depth information. The network can optimize depth information fitting for more precise prediction results. The Adam solver is used for training, with the default beta parameter of $(0.9, 0.999)$ and a learning rate of 5×10^{-4} .

We implemented our network in the PyTorch [PGM*19] framework and integrated the network into our distributed ray-tracer using the officially provided libtorch library. For the visibility and depth networks, the datasets are similar in size, although the content of the two datasets is different. We generate about 500 million ray samples for each base mesh. For different complexity geometries, the convergence speed of the two networks is also inconsistent, requiring approximately 3-5 hours using an NVIDIA RTX 3090 GPU. However, our neural proxy is scene-independent, so a trained model can be reused in different scenes.

5. Neural Proxy-based Ray Traversal

To reduce the ray transmission overhead in data-parallel ray tracing, we propose a neural proxy-based data-parallel ray traversal

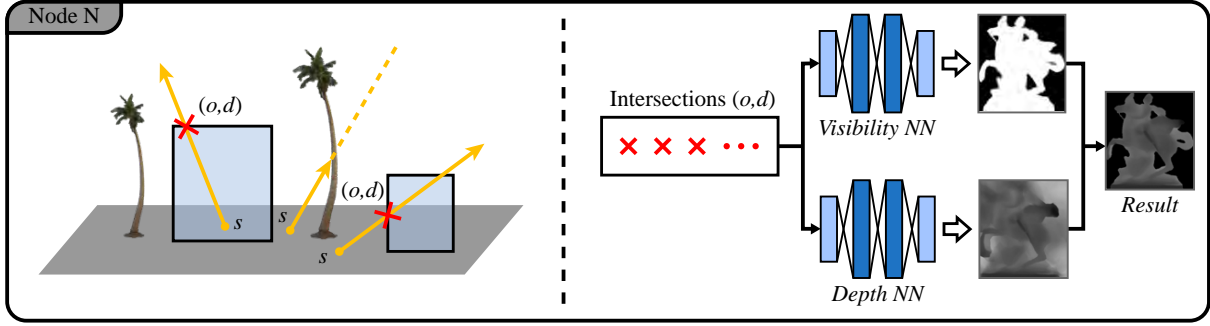


Figure 5: Proxy-based Ray Traversal. The left side shows the intersection of the rays with the node’s local geometry and the top-level BVH. On the right, the recorded intersection is fed into the neural proxy. The visibility NN and Depth NN correspond to the visibility neural network and the depth neural network, respectively.

method. In this approach, each node can use the cached neural proxies to determine the global intersection of rays without accessing the actual geometry. Consequently, a ray can be sent directly to the node containing the geometry of its closest intersection, significantly improving the transmission efficiency. Our neural proxy-based ray traversal method consists of two main stages: top-level BVH traversal and neural proxy intersection, as illustrated in Fig.5.

5.1. Top-level BVH Traversal

During top-level BVH traversal, the rays will traverse the top-level BVH at its generated node and record the intersection information. Initially, the rays intersect with the BVH of local objects and update their individual intersection distance t_{Max} and the index of the *next* geometry to be visited according to the intersection result. Subsequently, rays are intersected with the top-level BVH to obtain the non-local instance with which they intersect, along with the intersect positions o on bounding boxes and ray direction d , as illustrated in Fig.5. This gathered information will be used as input to the neural proxies to evaluate the global transmission targets of the rays.

To improve the processing efficiency, the intersection point o recorded for each ray must follow some rules. Firstly, for the same bounding box (leaf of the top-level BVH), we only record the closest intersection of a ray, as shown on the left side of Fig.6. Secondly, when the ray’s origin s lies inside the bounding box (Fig.6 (right)), we record the closest intersection o and an inversed direction d . In addition, the distance t from the ray origin s to o should also be recorded, denoted as t_{bbox} . Finally, it is crucial to note that to share the neural proxy among its related instances, all intersections must be inversely transformed to the base mesh’s local space. Accordingly, all data needs to be normalized.

5.2. Neural Proxy-based Ray Intersection

Following the top-level BVH traversal stage, the recorded intersections and directions (o, d) are input into the respective neural proxy to evaluate the visibility and depth of the actual geometries, as shown in Fig.5. Then, based on the visibility, depth, and top-level BVH intersection results, we can get the intersection distance

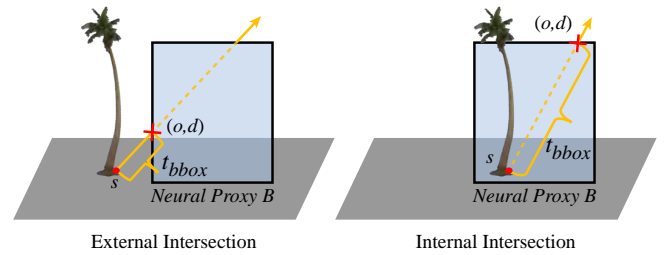


Figure 6: Illustration of external and internal intersections.

$t_{predicted}$ of the corresponding non-local geometry. Since the rays have already intersected with the node’s local geometries, we can obtain the rays’ intersection result with the entire scene once all the $t_{predicted}$ of non-local geometries have been evaluated.

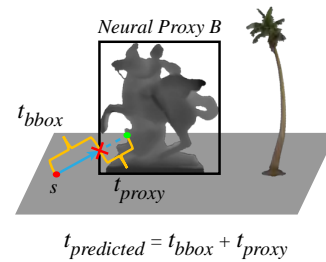


Figure 7: Intersection distance of external intersection.

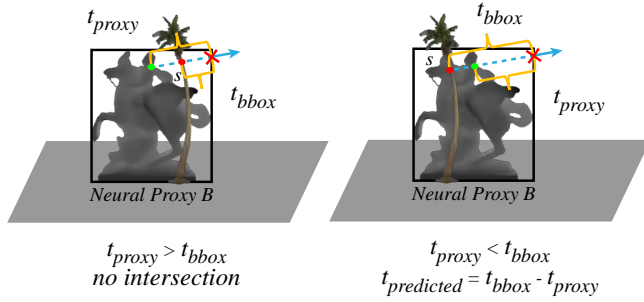
We use different methods to compute $t_{predicted}$ for the external and internal intersections respectively. For the external intersection, where the ray’s origin s lies outside the bounding box, as shown in Fig. 7. If the visibility value falls below a specified threshold, the ray is deemed not to intersect with the geometry represented by this neural proxy. Otherwise, the ray is determined to intersect the geometry, the *next* geometry to be visited and the intersection distance $t_{predicted}$ are updated. Specifically, the *closest-hits* is updated according to $t_{predicted}$, where $t_{predicted} = t_{bbox} + t_{proxy}$ with t_{proxy} representing the depth predicted by the neural proxy.

Table 1: Scene statistics of test scenes.

Scene	City	Sculpture	Airport	Courtyard	Moana
Instance	152	27	20	151	40
Instanced triangle (M)	1560	1970	151	57	13000
Base mesh	11	3	3	10	6
Texture (MB)	1679	194	810	134	904
Most GPU memory usage among nodes (GB)	11.93	7.24	10.94	8.96	8.74

* From top to bottom, the table presents the number of instances and instanced triangles, the number of base meshes, the size of the texture, and the most GPU memory usage among nodes.

As for the internal intersection, if the predicted visibility value is less than the specified threshold, it is conclusively determined that no intersection occurs. Instead, further evaluation is required using the t_{proxy} , as shown in Fig. 8. If t_{proxy} exceeds the t_{bbox} , then no intersection is indicated. Otherwise, the ray is considered to intersect the geometry, with the $t_{predicted} = t_{bbox} - t_{proxy}$. Then, the *next* geometry to be visited and the *closest-hits* are updated according to $t_{predicted}$.

**Figure 8:** Intersection distance of two different cases of internal intersection.

Our neural proxy-based ray traversal method can be used for primary, secondary, and shadow rays of data-parallel ray tracing. However, the introduction of our neural proxies will inevitably introduce errors in the rendering results. Additionally, owing to the multiple bounces of rays, secondary and shadow rays account for the major transmission overhead. Therefore, we chose to apply our proxy-based ray traversal for secondary and shadow rays. Thus, the secondary ray can determine which geometry its *closest-hits* is in and be forwarded directly to the node containing the geometry, thereby avoiding ineffective ray transmission. The shadow ray can directly determine occlusion and complete shadow computations at the node where they are generated without any ray transmission.

6. Result

In this section, we evaluate the performance of our proposed methods, including our entire data-parallel ray tracing framework and a detailed analysis of our neural proxy.

Our distributed framework is established on a commercial HPC cloud computing environment, where each node is equipped with an Intel SkyLake 6151 CPU and an NVIDIA Quadro RTX 5000

GPU with 16 GB device memory. The network bandwidth among nodes is 1Gig-E Ethernet, commonly used in personally built cluster conditions. Most of our tests were conducted using two distributed nodes. However, as shown in Table. 2, we employed more nodes to evaluate the scalability of our method.

We used five different types of scenes for our tests, and the statistics of the scenes are shown in Table. 1. For the scene Moana Island, its dominant storage overhead is from the terrain, which is unsuitable for our neural proxy. Consequently, we simplified the terrain and selected five typical objects to apply our neural proxy. Compared to the actual meshes, the memory footprint of our neural proxy is almost negligible, as shown in Fig. 9.

The results shown in this paper are all rendered using a 1920×1080 resolution, 512 samples per pixel (spp), and five maximum bounces. Each intersection point samples four shadow rays on the light source.

6.1. Performance and Quality Comparison

We use the state-of-the-art data-parallel ray tracing method proposed by Wald et al. [WP22] as the baseline for performance comparison and their rendering results as a reference for quality evaluation. It should be noted that the results of Wald et al. [WP22] have no quality loss compared to the ground truth. To evaluate image quality, we apply four metrics-peak signal-to-noise (PSNR), structural similarity (SSIM), FLIP [ANAM*20], and perceptual metric (LPIPS) [ZIE*18]-on the tone-mapped results produced by our

1936	1728	7198	5191	9088
1.19	1.19	1.19	1.19	1.19
3891	865	204	134	277
1.19	1.19	1.19	1.19	1.19

Figure 9: Storage comparison. The top value is the memory footprint (MB) of actual geometry and the bottom value is that of our neural proxy.

method and the reference. The rendering results and quality evaluation are illustrated in Fig.10 and Table.3. These results indicate that, owing to the advantages of our neural proxy in secondary and shadow ray transmission and computation, our method provides $2.29\sim 3.36\times$ speedup while maintaining a high rendering quality. It is evident that our neural proxy method can generate shadows

that are highly similar to the reference, with negligible quality loss in indirect lighting.

For the scalability comparison, we subdivided the base mesh of the City scene to generate scenes of varying scales. Then, different numbers of nodes were used to render these scenes separately, with the scene scales and rendering times shown in Table 2. It can be

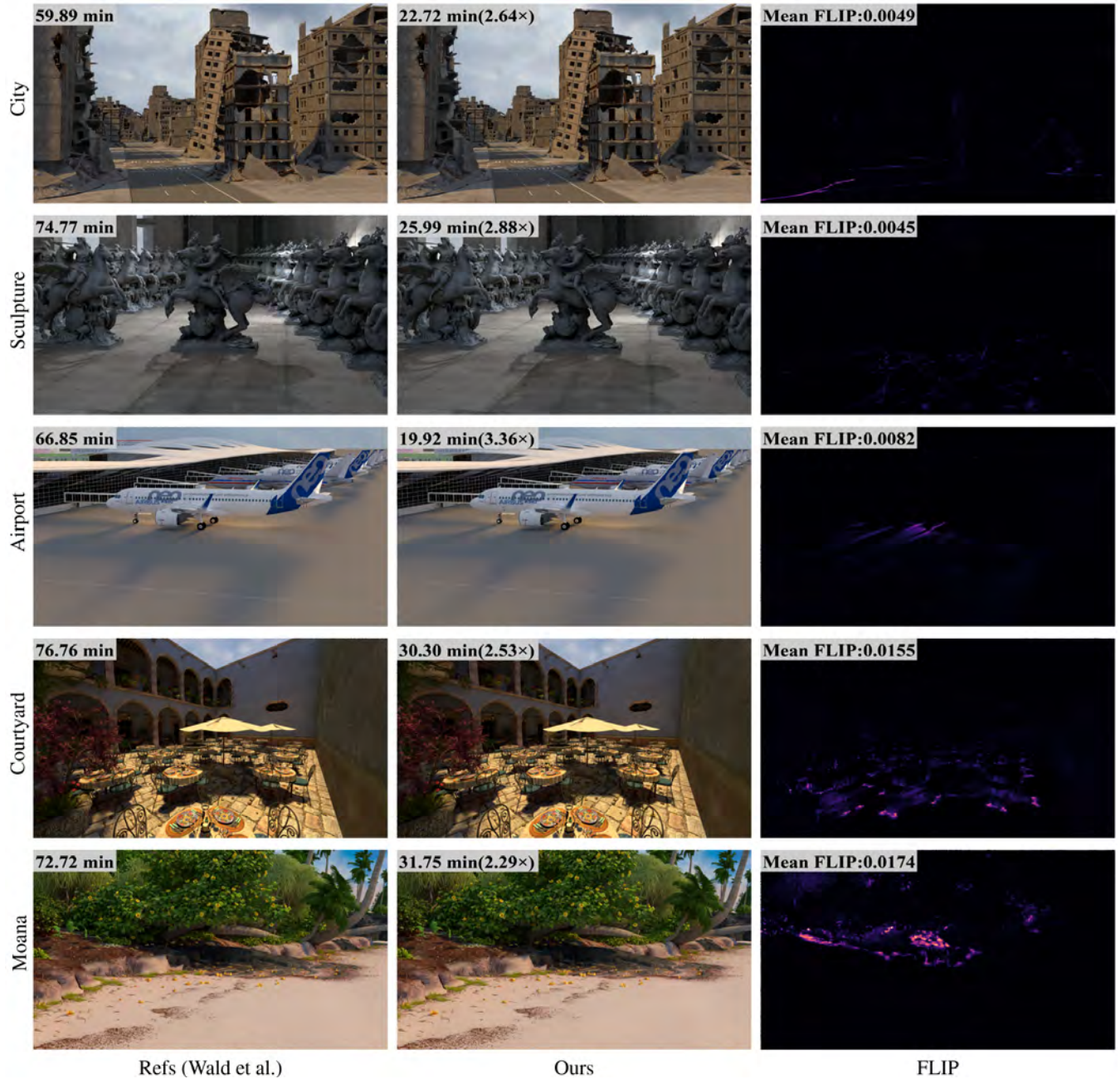


Figure 10: Performance and quality comparison. The images were rendered at 1920×1080 resolution, 512 spp, and 5 maximum bounces. The bandwidth among nodes is 1Gig-E.

seen that as the number of nodes increases, the rendering time of Wald et al. [WP22] increases dramatically, whereas the rendering time of our method is hardly affected by the number of rendering nodes. This is due to the fact that in Wald et al. [WP22] method, a ray typically needs to be transmitted multiple times between nodes to determine the *closest-hits*. Consequently, the number of times the ray is transmitted increases significantly with the number of nodes. In contrast, owing to the proxy-based ray traversal in our method, the secondary ray can be directly forwarded to the target node, regardless of the number of nodes, and the shadow ray does not need to be transmitted.

6.2. Performance Analysis at Different Bandwidths

Given that our approach’s efficacy relies on reducing ray transmission across nodes, we test the performance of our method under different network bandwidth environments, as shown in Table.4. It can be seen that smaller bandwidths lead to higher communication costs, while the benefits of our approach tend to be greater.

We also analyzed how our approach leads to speedups by measuring the time in each phase in Table.4. Undoubtedly, our method demonstrates a reduced data transmission volume to Wald et al. [WP22], resulting in a speedup range from $1.94\times$ to $3.78\times$ for data transmission time. As for ray intersections, despite the additional neural network evaluation time we introduce to our method, speedups ranging from $2.47\times$ to $2.49\times$ were achieved.

Table 2: Scalability comparison.

Scene scale(GB)	16.95	24.45	32.10	41.70	53.20
Node count	2	3	4	6	8
Rendering time(s)					
Ours	2.66	2.86	2.93	2.84	2.83
Wald et al.	7.02	9.12	10.83	13.58	15.12

The overhead caused by neural proxy and the speedup caused by avoiding the intersection of shadow rays and actual meshes is relatively constant. However, the speedup associated with transmission is related to the bandwidth among nodes. As a result, the speed of our method is inversely proportional to bandwidth.

6.3. Performance Analysis of Our Neural Proxy

We analyzed the storage efficiency as well as the accuracy of our neural proxy and the impact of its hyperparameters on the final rendering results. Fig. 11 shows the comparison with other geometric distribution representations, including classical edge collapse [GH97] and sparse voxel directed acyclic graphs (SVDAGs) [KSA13]. The reference is the unsimplified geometry, and we can see that our neural proxy has the best compression rate and accuracy. Although our neural proxy requires more training time, its training cost is acceptable, considering that the trained model can be used repeatedly.

The number of neurons per layer N is the hyperparameter that has the highest impact on the evaluation speed and the accuracy of our neural networks. Fig. 14, illustrates the impact of varying the sizes N of the visibility network on overall rendering performance and quality. As N increases, the rendering quality improves significantly, especially when increasing N from 64 to 128, but at the same time, the network evaluation overhead increases. The choice of network size requires a trade-off between quality and additional overheads. Eventually, we set the network size N to 256 for both the visibility network and the depth network.

The encoding method is also an important factor affecting the performance of neural implicit geometry. To explore this, we compared our direct encoding with the hash coding method [MESK22] used by Weier et al. [WRM*24], as shown in Fig. 12. It can be seen that the accuracy of the direct encoding method surpasses that of hash coding. This is attributed to the nature of our input data, which

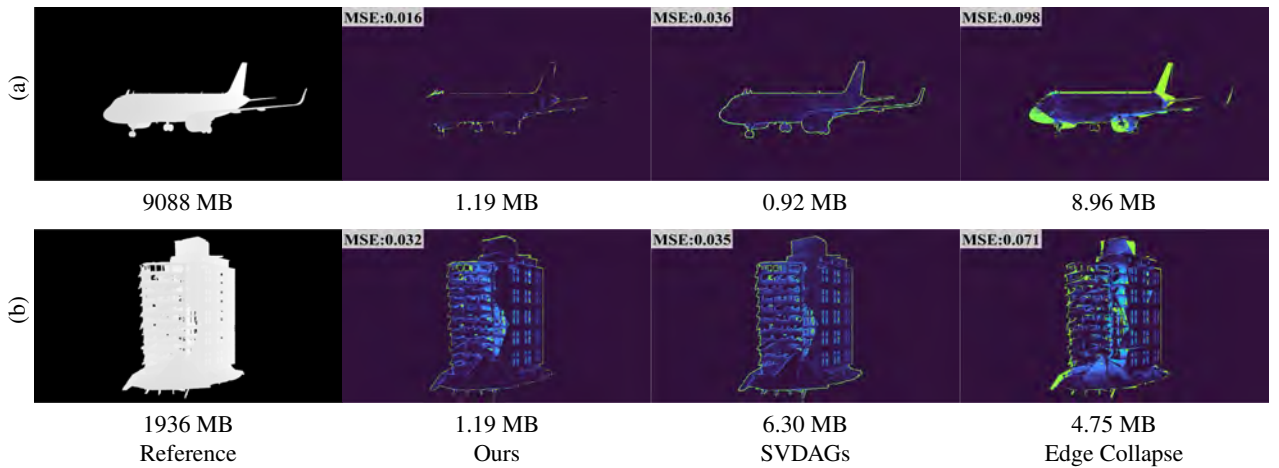


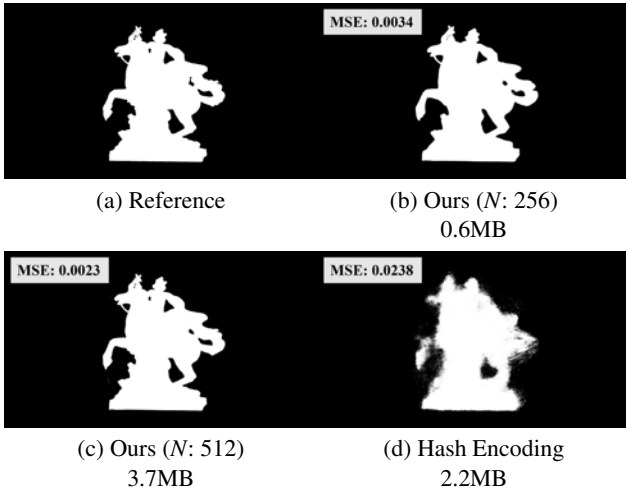
Figure 11: Comparison between different geometric distribution representations. The first column is utilized as the reference for visibility and depth, with the subsequent three columns showcasing the respective error between each method and the reference. Below each line is presented the storage required for each method.

Table 3: Quality comparison of five test scenes.

Scene	City	Sculpture	Airport	Courtyard	Moana
PSNR(dB)↑	49.08	54.87	54.77	45.46	47.04
SSIM↑	0.9995	0.9992	0.9991	0.9987	0.9979
Mean FLIP↓	0.0049	0.0045	0.0082	0.0155	0.0174
LPIPS↓	0.0007	0.0003	0.0002	0.0007	0.0026

Table 4: Time-consuming at 1spp in each phase of the Sculpture scene in different bandwidths. The fifth column contains the times for shadow ray processing, secondary ray processing, and intersection, respectively.

Scene	Bandwidth (Gig-E)	Method	Transmission time (s)	Neural Proxy & Intersection time (s)	Total time (s)
Sculpture	1	Our	0.927 (3.78×)	0.971 + 0.171 + 0.975 (2.48×)	3.044 (2.88×)
		Wald et al.	3.505	0 + 0 + 5.257	8.762
	2	Our	0.563 (3.36×)	0.967 + 0.161 + 0.980 (2.49×)	2.671 (2.68×)
		Wald et al.	1.889	0 + 0 + 5.257	7.146
	4	Our	0.431 (2.59×)	0.968 + 0.171 + 0.978 (2.48×)	2.548 (2.49×)
		Wald et al.	1.116	0 + 0 + 5.251	6.367
	20	Our	0.343 (1.94×)	0.967 + 0.171 + 0.983 (2.47×)	2.464 (2.39×)
		Wald et al.	0.666	0 + 0 + 5.233	5.903

**Figure 12:** Comparison of different encoding methods for visibility networks. (b) and (c) are our direct encoding methods with sizes N of 256 and 512, respectively. (d) is the result of the hash encoding method.

consists of positions on the bounding box and ray directions, rather than the sparse 3D positions that hash coding typically addresses.

6.4. Dynamic Scene

Our neural proxy is constructed in geometry’s local space and necessitates retraining only in the event of local space deformations. Consequently, it can generate stable and coherent rendering results of dynamic scenes with moving instances and light sources, as shown in Fig. 15. The visual coherence of dynamic scenes can also be evaluated in the provided video.

7. Limitations and Future Work

When the ray’s origin is located within a geometry’s bounding box, the neural proxy predicts the intersection depth $t_{predicted} = t_{bbox} + t_{proxy}$, as shown in Fig. 8. However, the ray might be occluded before reaching the predicted position, resulting in an incorrect depth estimation. While this issue does not impact our shadow generation, it may cause secondary rays to be sent to incorrect nodes. Additionally, our neural proxy currently struggles to capture high-frequency geometric details, leading to a noticeable quality loss in shadow, see Fig. 13.

In future work, we will subdivide and shrink the geometry’s bounding box, allowing for the construction of neural proxies on more compact bounding boxes [WRM*24]. This approach will minimize the empty space within the bounding box, thus avoiding internal intersection errors and improving the accuracy of the model. At the same time, we aim to expand the applicability of neural proxy to deal with primary ray traversal. Of course, this will require more complex neural networks to provide more and finer geometric information and will require us further to balance the benefits and overhead of the neural networks. Moreover, the adaptive selection of hyperparameters based on the complexity of the geometry is also an important research direction for improving the practicality of our method.

8. Conclusion

In this paper, we proposed a lightweight neural geometric representation used for massive scenes, termed a neural proxy, which can be shared among instances and utilized in dynamic scenes. Additionally, we proposed neural proxy-based data-parallel ray traversal methods that drastically reduced transmission and intersection overheads of secondary and shadow rays. Compared with the state-of-the-art, our neural proxy-based data-parallel ray tracing frame-



Figure 13: Rendering results of high-frequency shadow.

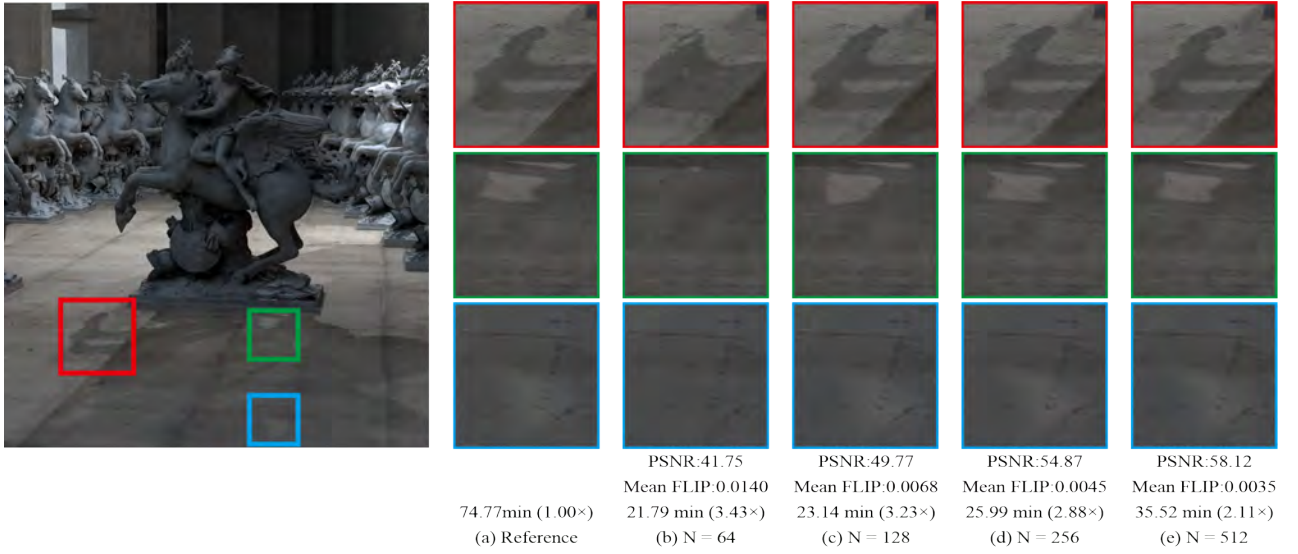


Figure 14: Comparison of different visibility network sizes N . The visibility network size N increases from left to right.

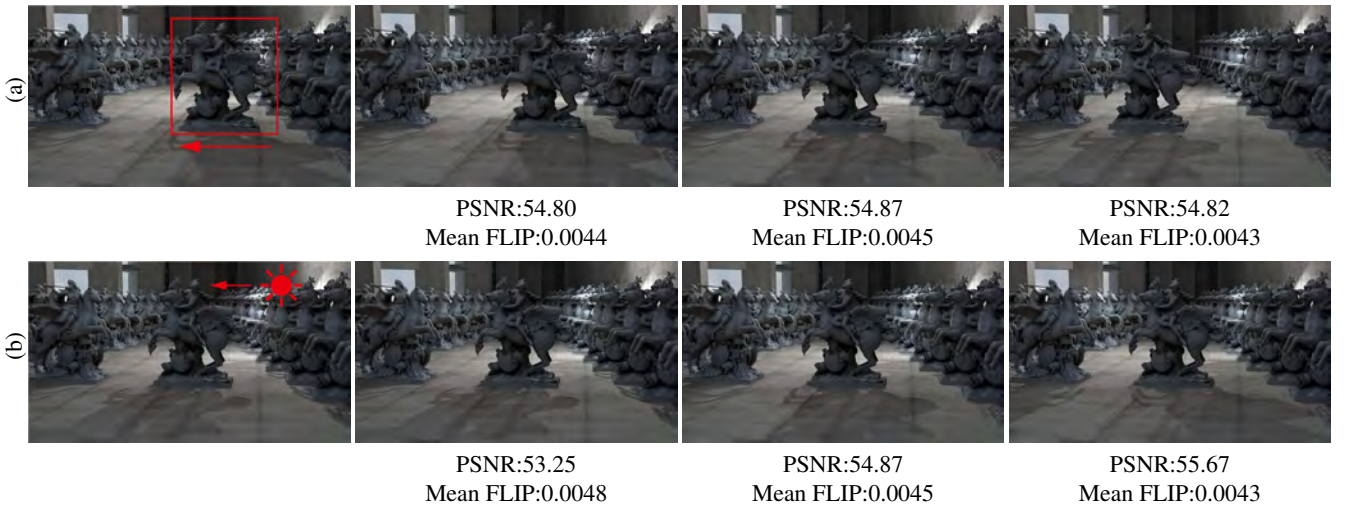


Figure 15: The result of dynamic scenes. In row (a), the position of the instance has moved from right to left. In row (b), the position of the light source has moved from right to left.

work yields a notable $2.29 \sim 3.36\times$ speedup while maintaining the quality of the rendered results.

Acknowledgements

We thank the reviewers for their valuable comments. This work has been partially supported by the National Key R&D Program of China under grant No.2022YFB3303200, the National Natural Science Foundation of China under grant No.62272275, the Shandong Provincial Natural Science Foundation under grant No.ZR2023QF123.

References

- [ANAM*20] ANDERSSON P., NILSSON J., AKENINE-MÖLLER T., OSKARSSON M., ÅSTRÖM K., FAIRCHILD M. D.: Flip: A difference evaluator for alternating images. *Proc. ACM Comput. Graph. Interact. Tech.* 3, 2 (2020), 15–1. [7](#)
- [Car05] CARUCCI F.: Inside geometry instancing. *GPU Gems 2* (2005), 47–67. [1, 2](#)
- [DGBP05] DEMARLE D. E., GRIBBLE C. P., BOULOS S., PARKER S. G.: Memory sharing for interactive ray tracing on clusters. *Parallel Computing* 31, 2 (2005), 221–242. [2](#)
- [FKH23] FUJIEDA S., KAO C. C., HARADA T.: Neural Intersection Function. In *High-Performance Graphics - Symposium Papers* (2023), Bikker J., Gribble C., (Eds.), The Eurographics Association. [doi:10.2312/hpg.20231135. 3](#)
- [FSP*22] FOULADI S., SHACKLETT B., POMS F., ARORA A., OZDEMIR A., RAGHAVAN D., HANRAHAN P., FATAHALIAN K., WINSTEIN K.: R2e2: Low-latency path tracing of terabyte-scale scenes using thousands of cloud cpus. *ACM Transactions on Graphics (TOG)* 41, 4 (2022), 1–12. [2](#)
- [GH97] GARLAND M., HECKBERT P. S.: Surface simplification using quadric error metrics. In *Proceedings of the 24th annual conference on Computer graphics and interactive techniques* (1997), pp. 209–216. [9](#)
- [GIF*18] GEORGIEV I., IZE T., FARNSWORTH M., MONTOYA-VOZMEKIANO R., KING A., LOMMEL B. V., JIMENEZ A., ANSON O., OGAKI S., JOHNSTON E., ET AL.: Arnold: A brute-force production path tracer. *ACM Transactions on Graphics (TOG)* 37, 3 (2018), 1–12. [3](#)
- [IC14] IVSON P., CELES W.: Instanced rendering of massive cad models using shape matching. In *2014 27th SIBGRAPI Conference on Graphics, Patterns and Images* (2014), IEEE, pp. 335–342. [3](#)
- [Joh13] JOHANSSON M.: Integrating occlusion culling and hardware instancing for efficient real-time rendering of building information models. In *International Conference on Computer Graphics Theory and Applications* (2013), vol. 2, SCITEPRESS, pp. 197–206. [3](#)
- [JŘSŠ21] JAROŠ M., ŘÍHA L., STRAKOŠ P., ŠPETKO M.: Gpu accelerated path tracing of massive scenes. *ACM Transactions on Graphics (TOG)* 40, 2 (2021), 1–17. [2](#)
- [KNE*01] KATO T., NISHIMURA H., ENDO T., MARUYAMA T., SAITO J., CHRISTENSEN P. H.: *Parallel Rendering and the Quest for Realism: The Kilauea Massively Parallel Ray Tracer*, pp. Tech. rep., IV–1–59, 2001. [2](#)
- [KSA13] KÄMPE V., SINTORN E., ASSARSSON U.: High resolution sparse voxel dags. *ACM Transactions on Graphics (TOG)* 32, 4 (2013), 1–13. [9](#)
- [KWR*17] KELLER A., WÄCHTER C., RAAB M., SEIBERT D., VAN ANTWERPEN D., KORNDÖRFER J., KETTNER L.: The iray light transport simulation and rendering system. *ACM SIGGRAPH 2017 Talks* (2017), 1–2. [2](#)
- [LGZL*20] LIU L., GU J., ZAW LIN K., CHUA T.-S., THEOBALT C.: Neural sparse voxel fields. *Advances in Neural Information Processing Systems* 33 (2020), 15651–15663. [3](#)
- [MBRS*21] MARTIN-BRUALLA R., RADWAN N., SAJJADI M. S., BARRON J. T., DOSOVITSKIY A., DUCKWORTH D.: Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2021), pp. 7210–7219. [3](#)
- [MESK22] MÜLLER T., EVANS A., SCHIED C., KELLER A.: Instant neural graphics primitives with a multiresolution hash encoding. *ACM Trans. Graph.* 41, 4 (jul 2022). URL: <https://doi.org/10.1145/3528223.3530127>, [doi:10.1145/3528223.3530127. 3, 9](#)
- [MST*21] MILDENHALL B., SRINIVASAN P. P., TANCIK M., BARRON J. T., RAMAMOORTHY R., NG R.: Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM* 65, 1 (2021), 99–106. [3](#)
- [NCFL13] NAVRÁTIL P. A., CHILDS H., FUSSELL D. S., LIN C.: Exploring the spectrum of dynamic scheduling algorithms for scalable distributed-memory ray tracing. *IEEE Transactions on Visualization and Computer Graphics* 20, 6 (2013), 893–906. [2](#)
- [NFLC12] NAVRÁTIL P. A., FUSSELL D. S., LIN C., CHILDS H.: Dynamic scheduling for large-scale distributed-memory ray tracing. In *Eurographics Symposium on Parallel Graphics and Visualization (EGPGV)* (2012), The Eurographics Association, pp. 61–70. [doi:10.2312/EGPGV/EGPGV12/061-070. 2](#)
- [PBD*10] PARKER S. G., BIGLER J., DIETRICH A., FRIEDRICH H., HOBEROCK J., LUEBKE D., MCALLISTER D., MCGUIRE M., MORLEY K., ROBISON A., ET AL.: Optix: a general purpose ray tracing engine. *Acm transactions on graphics (tog)* 29, 4 (2010), 1–13. [3, 5](#)
- [PGM*19] PASZKE A., GROSS S., MASSA F., LERER A., BRADBURY J., CHANAN G., KILLEEN T., LIN Z., GIMELSHEIN N., ANTIGA L., ET AL.: Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems* 32 (2019). [5](#)
- [PSL*98] PARKER S., SHIRLEY P., LIVNAT Y., HANSEN C., SLOAN P.-P.: Interactive ray tracing for isosurface rendering. In *Proceedings Visualization'98 (Cat. No. 98CB36276)* (1998), IEEE, pp. 233–238. [2](#)
- [RCJ99] REINHARD E., CHALMERS A., JANSEN F. W.: Hybrid scheduling for parallel rendering using coherent ray tasks. In *IEEE Symposium on Parallel Visualization and Graphics (EGPGV)* (1999), IEEE, pp. 21–28. [2](#)
- [SKD*13] SLOMP M., KAWASAKI H., DARIBO I., SAGAWA R., FURUKAWA R., HIURA S., ASADA N.: Hardware-accelerated geometry instancing for surfel and voxel rendering of scanned 4d media. In *11th International Conference on Quality Control by Artificial Vision (QCAV)* (2013), vol. 3. [3](#)
- [SMB*20] SITZMANN V., MARTEL J., BERGMAN A., LINDELL D., WETZSTEIN G.: Implicit neural representations with periodic activation functions. *Advances in neural information processing systems* 33 (2020), 7462–7473. [3](#)
- [Stu18] STUDIOS W. D. A.: Moana island scene (v1.1) [data set]. <https://disneyanimation.com/data-sets/?drawer=/resources/moana-island-scene/>, 2018. [1](#)
- [SY17] SON M., YOON S.-E.: Timeline scheduling for out-of-core ray batching. In *Proceedings of High-Performance Graphics* (2017), pp. 1–10. [2](#)
- [TLY*21] TAKIKAWA T., LITALIEN J., YIN K., KREIS K., LOOP C., NOWROUZEZAHRAI D., JACOBSON A., MCGUIRE M., FIDLER S.: Neural geometric level of detail: Real-time rendering with implicit 3d shapes. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (Los Alamitos, CA, USA, jun 2021), IEEE Computer Society, pp. 11353–11362. URL: <https://doi.ieeecomputersociety.org/10.1109/CVPR46437.2021.01120>, [doi:10.1109/CVPR46437.2021.01120. 3](#)

- [Wal05] WALD I.: Realtime ray tracing and interactive global illumination. *Ausgezeichnete Informatikdissertationen 2004* (2005). 3
- [WBS03] WALD I., BENTHIN C., SLUSALLEK P.: Distributed interactive ray tracing of dynamic scenes. In *IEEE Symposium on Parallel and Large-Data Visualization and Graphics, 2003. PVG 2003.* (2003), IEEE, pp. 77–85. 2
- [WJZ23] WALD I., JAROŠ M., ZELLMANN S.: Data parallel multi-gpu path tracing using ray queue cycling. In *Computer Graphics Forum* (2023), Wiley Online Library, p. e14873. 2
- [WP22] WALD I., PARKER S. G.: Data parallel path tracing with object hierarchies. *Proceedings of High Performance Graphics* (2022). 1, 2, 4, 7, 9
- [WRM*24] WEIER P., RATH A., MICHEL E., GEORGIEV I., SLUSALLEK P., BOUBEKEUR T.: N-bvh: Neural ray queries with bounding volume hierarchies. In *ACM SIGGRAPH 2024 Conference Papers* (New York, NY, USA, 2024), SIGGRAPH '24, Association for Computing Machinery. URL: <https://doi.org/10.1145/3641519.3657464>, doi:10.1145/3641519.3657464. 3, 9, 10
- [WVB*14] WALD I., WOOP S., BENTHIN C., JOHNSON G. S., ERNST M.: Embree: a kernel framework for efficient cpu ray tracing. *ACM Transactions on Graphics (TOG)* 33, 4 (2014), 1–8. 3
- [XWPG*24] XU X., WANG L., PÉRARD-GAYOT A., MEMBARTH R., LI C., YANG C., SLUSALLEK P.: Temporal coherence-based distributed ray tracing of massive scenes. *IEEE Transactions on Visualization and Computer Graphics* 30 (2024), 1489–1501. 2
- [YKM*20] YARIV L., KASTEN Y., MORAN D., GALUN M., ATZMON M., RONEN B., LIPMAN Y.: Multiview neural surface reconstruction by disentangling geometry and appearance. *Advances in Neural Information Processing Systems* 33 (2020), 2492–2502. 3
- [ZBX*21] ZHU J., BAI Y., XU Z., BAKO S., VELÁZQUEZ-ARMENDÁRIZ E., WANG L., SEN P., HASAN M., YAN L.-Q.: Neural complex luminaires: representation and rendering. *ACM Trans. Graph.* 40, 4 (2021), 1–12. 5
- [ZIE*18] ZHANG R., ISOLA P., EFROS A. A., SHECHTMAN E., WANG O.: The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (2018), pp. 586–595. 7
- [ZWB*22] ZELLMANN S., WALD I., BARBOSA J., DERMIC S., SAHISTAN A., GÜDÜKBAY U.: Hybrid image-/data-parallel rendering using island parallelism. In *2022 IEEE 12th Symposium on Large Data Analysis and Visualization (LDAV)* (2022), IEEE, pp. 1–10. 2