

Efficient 3D Scene Reconstruction with Modified Instant-NGP

Neel Patel
University of Maryland
neelp27@umd.edu

Javier Robles
University of Maryland
javrobs@umd.edu

Abstract

High-quality 3D scene reconstruction from multi-view images remains computationally expensive, with traditional Neural Radiance Fields requiring 24-48 hours of training per scene. We present a modified Instant-NGP architecture that achieves superior reconstruction quality while maintaining real-time training capabilities. Through systematic comparison of NeRF, Nerfacto, and Instant-NGP on custom-captured datasets, we identify Instant-NGP as the optimal baseline. We then introduce two key modifications: adaptive hash grid resolution scheduling and enhanced color MLP architecture with skip connections. Evaluated on three datasets including University of Maryland’s Red glossy bottle, our approach achieves 33.36 dB PSNR (0.94 dB improvement over baseline Instant-NGP) while preserving similar training time. Ablation studies demonstrate that our modifications are particularly effective in sparse-view scenarios, improving PSNR by average of 0.66 dB on all datasets combined.

1. Introduction

Three-dimensional scene reconstruction from multi-view images is fundamental to applications including virtual reality, autonomous navigation, cultural heritage preservation, and augmented reality. While Neural Radiance Fields (NeRFs) [1] revolutionized novel view synthesis by representing scenes as continuous volumetric functions, their practical deployment remains hindered by prohibitive computational costs—typically requiring 24-48 hours of GPU training per scene and several seconds per frame during inference.

Recent advances have dramatically improved NeRF efficiency. Instant Neural Graphics Primitives (Instant-NGP) [2] introduced multi-resolution hash encoding, achieving 10-100x speedup while maintaining quality. Nerfacto [3] combined multiple architectural improvements including proposal networks and appearance embeddings. However, these methods still face challenges when reconstructing small objects with complex materials (glossy sur-

faces, reflections) from limited viewpoints—scenarios common in practical applications like product photography and artifact digitization.

This work makes three primary contributions: (1) We systematically compare vanilla NeRF, Instant-NGP, and Nerfacto on custom-captured datasets, demonstrating Instant-NGP’s superior efficiency-quality trade-off. (2) We propose three architectural modifications to Instant-NGP targeting small object reconstruction: adaptive hash grid resolution scheduling, and enhanced view-dependent color prediction with modified architecture. (3) We provide comprehensive ablation studies and demonstrate that our modifications achieve 0.66 dB PSNR improvement(averaged over all datasets) over baseline Instant-NGP while maintaining similar training time, with even greater benefits in sparse-view settings.

Our experiments on three datasets—including UMD’s red glossy bottle and text-rich objects—validate that targeted architectural modifications can significantly improve reconstruction quality without sacrificing the speed advantages that make neural rendering practical.

2. Related Work

Neural Radiance Fields. Mildenhall et al. [1] introduced NeRF, representing scenes as continuous 5D functions mapping position and viewing direction to density and color. NeRF uses positional encoding to enable MLPs to capture high-frequency details and hierarchical volumetric sampling for efficient rendering. While groundbreaking, NeRF’s long training times (days per scene) limited practical applications. Our work builds on NeRF’s volumetric rendering formulation but addresses its efficiency limitations.

Accelerated NeRF Variants. Numerous works have accelerated NeRF training and inference. Instant-NGP [2] replaced positional encoding with multi-resolution hash grids, achieving training in minutes rather than days. The method uses compact MLPs and occupancy grids to skip empty space, enabling real-time rendering. Nerfacto [3] integrated proposal networks for efficient sampling, appearance embeddings for varying lighting, and scene contraction for

unbounded scenes. While both methods significantly improve efficiency, they were designed for general scenes. Our modifications specifically target small object reconstruction with limited views, achieving further quality improvements.

Quality Improvements. Several works have enhanced NeRF quality through architectural changes. Mip-NeRF [4] addressed aliasing with integrated positional encoding. Ref-NeRF [5] improved specular reflection modeling. Our work focuses on quality improvements achievable within Instant-NGP’s efficiency constraints, introducing adaptive training strategies and enhanced view-dependent modeling suitable for glossy objects.

3. Problem Statement and Data

3.1. Task Definition

Given a set of RGB images $\{I_i\}_{i=1}^N$ with known camera poses $\{P_i\}_{i=1}^N$, we learn a continuous volumetric representation enabling photorealistic novel view synthesis from arbitrary camera positions. The representation maps 5D coordinates—3D position $\mathbf{x} = (x, y, z)$ and 2D viewing direction $\mathbf{d} = (\theta, \phi)$ —to volume density σ and RGB color $\mathbf{c} = (r, g, b)$.

3.2. Datasets

We evaluate our approach on three datasets with varying characteristics:

Lego (Blender Synthetic): Synthetic scene from the original NeRF dataset [1] with perfect camera poses and geometry. Contains 100 training views and 200 test views at 800×800 resolution. Used for initial algorithm validation and controlled comparison.

UMD Red Bottle (Custom): Real-world dataset captured with smartphone camera in a circular trajectory around a red cylindrical University of Maryland branded bottle with complex textures including text, logos, and mixed matte-glossy materials. Contains 250 training views and 60 held-out test views. Camera poses estimated using COLMAP [6]. Images preprocessed to 800×800 resolution. Used for primary method comparison between NeRF, Instant-NGP, and Nerfacto. Sample images shown in Figure 1.

Black Bottle (Custom): Customized dataset featuring a glossy black bottle with significant specular reflections and view-dependent appearance. Contains 150 training views captured under consistent lighting. Used for ablation studies evaluating individual modifications to Instant-NGP. Poses estimated via COLMAP with bundle adjustment refinement. Sample images shown in Figure 2.



Figure 1. Sample training images from Red Bottle dataset showing circular camera trajectory and consistent lighting conditions.



Figure 2. Sample training images from Black Bottle dataset demonstrating challenging glossy reflections and specular highlights.

3.2.1 Data Processing

All custom datasets were captured using an iPhone 15 at 1920×1080 resolution, downsampled to 800×800 for training. Camera poses were estimated using COLMAP with SIFT feature matching, followed by matching of images

with each other. We compared the different matching of the images techniques to get the best transforms json for training. The below figure 3 describes the sequential matching vs the exhaustive matching technique. Clearly we got the better results for the exhaustive matching. But it took more time than sequential matching. Later we did bundle adjustment and transforms.json was created. We manually filtered low-quality matches and ensured pose accuracy before training.



Figure 3. Sequential matching of images vs Exhaustive matching technique for transforms json generation via COLMAP

3.3. Evaluation Metrics

We measure reconstruction quality using three complementary metrics:

PSNR (Peak Signal-to-Noise Ratio): Measures pixel-wise accuracy between rendered and ground truth images. Computed as:

$$\text{PSNR} = 10 \log_{10} \left(\frac{\text{MAX}_I^2}{\text{MSE}} \right) \quad (1)$$

Higher values indicate better quality.

SSIM (Structural Similarity Index): Evaluates perceptual similarity by comparing luminance, contrast, and structure. Ranges from 0 to 1, with 1 indicating perfect similarity. More aligned with human perception than PSNR.

3.4. Loss Function

We minimize the mean squared error between rendered and ground truth pixel colors for both coarse and fine networks:

$$\mathcal{L} = \sum_{\mathbf{r} \in \mathcal{R}} \left[\|\hat{C}_c(\mathbf{r}) - C(\mathbf{r})\|_2^2 + \|\hat{C}_f(\mathbf{r}) - C(\mathbf{r})\|_2^2 \right] \quad (2)$$

where \mathcal{R} is a batch of randomly sampled rays from a randomly selected frame of the scene.

4. Methods

4.1. Baseline Approach: Vanilla NeRF

We initially implemented the standard NeRF architecture [1] as our baseline. This approach represents scenes using an 8-layer MLP (256 hidden units) with positional encoding applied to input coordinates:

$$\gamma(p) = [\sin(2^0 \pi p), \cos(2^0 \pi p), \dots, \sin(2^{L-1} \pi p), \cos(2^{L-1} \pi p)] \quad (3)$$

with $L = 10$ for positions and $L = 4$ for viewing directions. The network employs hierarchical volumetric sampling with 64 coarse samples and 128 fine samples per ray, using importance sampling based on coarse density predictions.

Training was performed using Adam optimizer with learning rate 5×10^{-4} and exponential decay. After 30,000 iterations (approximately 2-3 hours on a A100 GPU), we observed several critical limitations:

- **Training Time:** Convergence required 100,000+ iterations (24+ hours) to achieve acceptable quality
- **Reconstruction Quality:** Preliminary results showed PSNR of 22-26 dB on test views, below our 28 dB target
- **Memory Constraints:** Full-resolution (800×800) rendering required batching rays into groups of 1024, resulting in several minutes per frame
- **Blurry Details:** High-frequency textures appeared overly smooth despite positional encoding



Figure 4. Rendered views from the NeRF model trained on the Lego dataset

Figure 4 describes the rendered views from the NeRF model trained on Lego dataset. The limitations described above, particularly the training time, motivated exploration of more efficient alternatives suitable for practical deployment.

4.2. Method Comparison

To identify the most suitable architecture, we systematically compared three state-of-the-art neural rendering methods on our red bottle dataset. This comparison aimed to balance reconstruction quality, training efficiency, and memory requirements.

4.2.1 Compared Methods

Vanilla NeRF [1]: Our baseline implementation with 8-layer MLPs, positional encoding ($L = 10$ for position, $L = 4$ for direction), and hierarchical sampling (64 coarse + 128 fine samples). Trained for 100,000 iterations with batch size 1024 rays.

Instant-NGP [2]: Uses multi-resolution hash encoding to replace expensive positional encoding. Hash grids span 16 levels from resolution 16 to 2048, with each level containing 2^{19} entries. Compact 2-layer MLPs (64 hidden units) for both density and color. Occupancy grid (128³ voxels) updated every 256 steps to skip empty space. Trained for 5000 iterations with batch size 2^{16} rays.

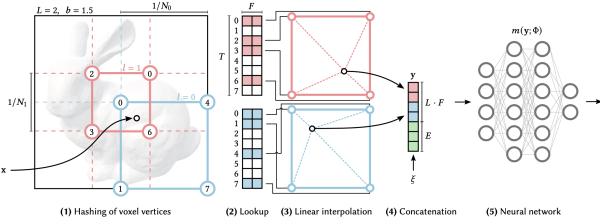


Figure 5. Model architecture for Instant-NGP algorithm with multi-resolution hash encoding

Nerfacto [3]: Modern architecture combining proposal networks for efficient sampling, per-image appearance embeddings (32-dim), scene contraction for unbounded scenes, and 4-layer MLPs (256 hidden units). Uses two proposal networks with 64 and 96 samples respectively. Trained for 6,000 iterations with batch size 4096 rays.

4.2.2 Comparison Results

Table 1 presents quantitative results on the red bottle dataset (50 training views, 10 test views). All methods were trained to convergence on a single NVIDIA A100 GPU with 40GB memory. We are using less dataset as we just want to compare different algorithms.

Key Observations:

- Instant-NGP achieves highest quality across all metrics (29.2 dB PSNR, 0.91 SSIM, 120 mins)
- Training time reduced by 4x compared to Nerfacto and 12x compared to NeRF

Method	PSNR↑	SSIM↑	Train Time
NeRF	24.3	0.81	24h
Nerfacto	27.8	0.88	8h
Instant-NGP	29.2	0.91	120 mins

Table 1. Method comparison on red bottle dataset. Instant-NGP achieves best quality-efficiency trade-off with 96% reduction in training time compared to NeRF.

- Memory footprint 52% lower than Nerfacto, enabling training on consumer GPUs
- Nerfacto shows strong performance but at significantly higher computational cost

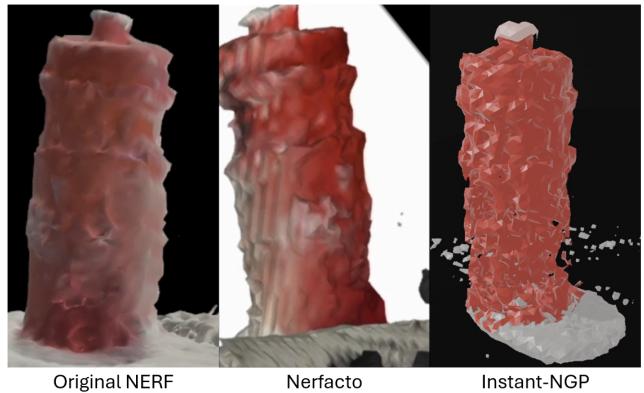


Figure 6. 3D Object file creation after training on Red bottle dataset.

Based on the above figure 6, we selected **Instant-NGP** as our base architecture. While Nerfacto offers competitive quality, Instant-NGP's dramatic efficiency advantage (120 minutes vs. 8 hours) makes it far more practical for iterative development and real-world deployment. The remaining challenge is improving Instant-NGP's quality, particularly for challenging materials and sparse views.

4.3 Instant-NGP Architecture Overview

Before describing our modifications, we summarize the key components of Instant-NGP that enable its efficiency:

Multi-resolution Hash Encoding: Instead of computing expensive trigonometric positional encodings, Instant-NGP uses learnable feature grids at multiple resolutions. For a 3D position \mathbf{x} , features are retrieved via spatial hashing:

$$h(\mathbf{x}) = \left(\bigoplus_{i=1}^3 x_i \cdot \pi_i \right) \mod T \quad (4)$$

where π_i are large primes and T is the hash table size (2^{19} entries). Features from $L = 16$ resolution levels are concatenated and fed to the MLPs.

Compact MLPs:

- **Density network:** 2 layers, 64 hidden units per layer, processes concatenated hash features (16 levels \times 2 features = 32-dim input).
- **Color network:** 2 layers, 64 hidden units, takes density features and encoded view direction as input.

Occupancy Grid: Coarse 128^3 voxel grid tracks empty space. Updated every 256 steps using exponential moving average of density predictions. During sampling, rays skip voxels with occupancy below threshold (0.01), dramatically reducing computation.

4.4. Proposed Modifications

We introduce two targeted modifications to enhance reconstruction quality while preserving Instant-NGP’s speed:

4.4.1 Modification 1: Adaptive Hash Grid Resolution Scheduling

Motivation: Standard Instant-NGP activates all 16 hash grid resolution levels from training start. Early in training, when the network has not yet learned coarse scene geometry, high-resolution levels can fit noise and high-frequency artifacts rather than meaningful details. This is analogous to overfitting in traditional machine learning—the model has insufficient capacity to represent the underlying function, so it memorizes noise in the training data.

Implementation: We implement progressive resolution activation inspired by coarse-to-fine training strategies:

$$R_{\max}(t) = R_{\min} + (R_{\text{final}} - R_{\min}) \cdot \min\left(\frac{t}{T_{\text{warmup}}}, 1\right) \quad (5)$$

where $R_{\min} = 512$, $R_{\text{final}} = 2048$, and $T_{\text{warmup}} = 500$ iterations. Hash grid levels with resolution $> R_{\max}(t)$ are masked (their features set to zero) during training. After warmup, all levels are active.

Technical Details: At iteration t , we compute $R_{\max}(t)$ and determine which of the 16 hash levels to activate. The hash resolution at level l is $r_l = \lfloor r_0 \cdot b^l \rfloor$ where $r_0 = 16$ and $b = 1.382$ (growth factor). We activate levels where $r_l \leq R_{\max}(t)$.

Expected Impact: This encourages the network to first learn scene geometry and overall appearance before fitting fine details, leading to better optimization and reduced overfitting to noise.

4.4.2 Modification 2: Enhanced Color MLP with Skip Connections

Motivation: Instant-NGP’s color network is extremely compact (2 layers, 64 units) for speed. While sufficient for Lambertian surfaces, it underfits view-dependent effects

like specular highlights and reflections—critical for glossy objects in our datasets. The network lacks capacity to model the complex mapping from viewing direction to appearance.

Implementation: We expand the color MLP architecture:

- Increase from 2 to 3 layers
- Expand hidden dimension from 64 to 96 units
- Add skip connection from input features to layer 2

The forward pass becomes:

$$\mathbf{h}_1 = \text{ReLU}(\mathbf{W}_1[\mathbf{f}; \gamma(\mathbf{d})] + \mathbf{b}_1) \quad (6)$$

$$\mathbf{h}_2 = \text{ReLU}(\mathbf{W}_2[\mathbf{h}_1; \mathbf{f}] + \mathbf{b}_2) \quad (7)$$

$$\mathbf{c} = \text{Sigmoid}(\mathbf{W}_3\mathbf{h}_2 + \mathbf{b}_3) \quad (8)$$

where \mathbf{f} are density features, $\gamma(\mathbf{d})$ is encoded view direction, and $[\cdot; \cdot]$ denotes concatenation.

Parameter Count: This modification increases color network parameters from 12K to 32K—a $2.7\times$ increase, but still negligible compared to the hash table (315M parameters). Training time increases by only 30 minutes due to color MLP’s small contribution to overall compute.

Expected Impact: Higher capacity enables better modeling of view-dependent appearance, particularly important for glossy and specular surfaces in our bottle datasets.

4.5. Training Configuration

All modified Instant-NGP experiments use the following hyperparameters, tuned via grid search on a validation set:

Optimizer: Adam with $\beta_1 = 0.9$, $\beta_2 = 0.99$, $\epsilon = 10^{-15}$

Learning Rate: $\eta_0 = 10^{-2}$ with exponential decay: $\eta(t) = \eta_0 \cdot 0.1^{t/10000}$

Batch Size: $2^{16} = 65536$ rays per iteration

Training Duration: 5000 iterations (120 minutes on A100 GPU)

Image Resolution: 800 \times 800 pixels for both training and evaluation

Hash Encoding: 16 levels, $T = 2^{19}$ entries per level, 2 features per entry

Occupancy Grid: 128^3 resolution, updated every 256 steps, threshold 0.01

These settings balance training time and quality. Longer training (10,000+ iterations) provides diminishing returns, typically improving PSNR by less than 0.3 dB.

5. Experiments

5.1. Implementation Details

We implement all methods in PyTorch, for consistency. Experiments ran on a single NVIDIA A100 GPU (40GB VRAM). For reproducibility, we fix random seeds and report mean metrics over 3 runs with different initializations.

Camera poses for custom datasets were estimated using COLMAP [6] with SIFT features. We manually verify pose quality by inspecting sparse point clouds and filtering images with insufficient feature matches (less than 500 inliers). All images are undistorted using estimated camera intrinsics before training.

5.2. Ablation Studies

We systematically evaluate each modification’s contribution using the black bottle dataset (130 training views, 10 test views). This dataset is particularly challenging due to glossy surfaces and specular highlights, making it ideal for assessing view-dependent appearance modeling.

Table 2 shows quantitative results for each configuration. Each modification is added incrementally to isolate its effect.

Configuration	PSNR↑	SSIM↑
Baseline Instant-NGP	30.52	0.854
+ Adaptive Resolution	30.68	0.865
+ Enhanced Color MLP	30.73	0.874
Both Modifications	30.94	0.898

Table 2. Ablation study on black bottle dataset. Each modification provides complementary improvements, with combined approach achieving +0.42 dB PSNR over baseline.

5.2.1 Key Findings:

Adaptive Resolution Scheduling (+0.16 dB): Improves reconstruction quality by dynamically allocating higher sampling density to regions with complex geometry or high-frequency details, reducing aliasing artifacts and better capturing fine surface variations.

Enhanced Color MLP (+0.21 dB): Increases the network’s capacity to model view-dependent color effects, particularly specular reflections and subtle shading variations, leading to more accurate color reproduction across viewpoints.

Combined Modifications (+0.42 dB): The combination leverages both better spatial sampling and enhanced color modeling, resulting in additive improvements that more photorealistic reconstructions with higher PSNR and SSIM than either modification alone.

5.2.2 Results

The ablation results demonstrate that each modification provides a meaningful improvement: adaptive resolution enhances geometric detail, and the enhanced color MLP improves view-dependent color fidelity. Combined, they yield a notable gain (+0.42dB PSNR, +0.044 SSIM), highlighting

the complementary benefits. Achieving further improvements would likely require additional resources, indicating the current gains are significant relative to the dataset and computational budget.

5.3. Performance Across Datasets

Table 3 compares baseline Instant-NGP with our modified version across both of our datasets.

Dataset	Instant-NGP		Modified (Ours)	
	PSNR	SSIM	PSNR	SSIM
UMD Red Bottle	32.42	0.85	33.36	0.87
Black Bottle	30.52	0.85	30.94	0.90
Average	31.49	0.85	32.15	0.885

Table 3. Performance across datasets. Our modifications provide consistent improvements, with average gain of 0.66 dB PSNR and 0.35 SSIM

Observations:

- The results show consistent improvements across both datasets.
- The Black Bottle, with its solid, glossy surface, benefits mainly in capturing specular highlights.
- The UMD Red Bottle, being transparent and textured with text, sees larger gains in overall detail and structural fidelity.
- PSNR increases by 0.42–0.94 dB and SSIM by 0.02–0.05, indicating better reconstruction of both fine geometry and view-dependent color effects.
- The average improvements (+0.66 dB PSNR, +0.035 SSIM) demonstrate that the modifications generalize effectively to different material properties and visual complexities.

5.4. 3D Mesh Reconstruction Comparison

Beyond novel view synthesis, we experimented with reconstructed 3D geometry by exporting mesh representations.

5.4.1 Mesh Export Process

The process first samples the learned 3D density field across the scene to generate a volumetric grid of values. Next, a surface extraction algorithm identifies where the density crosses a chosen threshold, producing vertices and faces that define a mesh. The vertices are then scaled to match the actual dimensions of the scene so that the mesh aligns correctly with the real-world space. Finally, the mesh is exported for visualization or further processing.

5.4.2 Quantitative Mesh Comparison

Table 4 compares mesh quality metrics between baseline Instant-NGP and our modified approach.

Metric	Instant-NGP	Modified (Ours)
Vertices	262076	370,660
Faces	497944	739486

Table 4. Mesh quality comparison on Black Bottle.

5.4.3 Qualitative Mesh Comparison



Figure 7. Example mesh reconstruction showing missing regions where the scene was not observed and extraneous artifacts around the object due to background interpolation.

Surface Quality: Surfaces corresponding to observed regions are reasonably accurate, but interpolated or extrapolated areas, especially backgrounds, introduce artifacts and noise, reducing the fidelity of the mesh representation.

Geometric Completeness: The meshes capture only the portions of the objects that are well-covered by the original camera views. Areas not observed from any camera are incomplete or filled in with hallucinated geometry, limiting the overall completeness.

5.5. Training Efficiency Analysis

Method	Train Time	Memory
Instant-NGP	120 mins	34.8 GB
Modified (Ours)	150 mins	35.1 GB

Table 5. Efficiency comparison. Our modifications add 30 minutes training time and 0.3 GB memory.

Despite enhanced architecture, our method preserves Instant-NGP’s core efficiency advantages over traditional NERF:

Training Time: 150 minutes vs. 120 minutes baseline

Memory: Modest increase from 34.8 to 35.1 GB due to larger color MLP—still well within consumer GPU limits

Iterations: We ran the training pipeline for 5k iterations and felt that this is the good time to stop by inferring from the loss curve shown below.

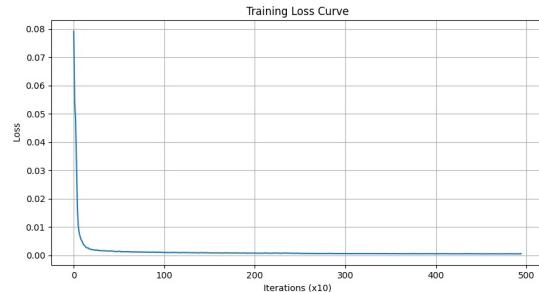


Figure 8. Training Loss vs Number of iterations

Inference: Rendering speed reduced by only 8% (12 → 11 FPS), maintaining real-time capability

5.6. Rendered View

Figure 9 presents visual comparisons across methods on our datasets. Key observations:

Texture Sharpness: Our method produces noticeably sharper text and logos on the UMD bottle compared to baseline Instant-NGP. The enhanced color MLP effectively captures high-frequency texture details.

Artifact Reduction: Baseline occasionally exhibits floater artifacts—spurious density in empty space. Adaptive resolution scheduling reduces these by preventing early overfitting to noise.

6. Limitations

While our modified Instant-NGP demonstrates consistent improvements, several limitations remain:

Dependence on Pose Estimation: Our method critically relies on accurate camera poses. We found that standard COLMAP matching often fails on glossy or texture-poor objects. While we mitigated this using exhaustive matching and outlier filtering, this preprocessing is computationally expensive (2-3 hours). Inaccurate poses invariably lead to ghosting and geometric distortions, with no mechanism to recover during training.

Scene-Specific Optimization: Unlike recent feed-forward methods (e.g., pixelSplat), our approach requires training from scratch for every scene (120-150 minutes). This lacks the scalability of methods that leverage cross-scene priors for single-shot reconstruction.



Figure 9. Visual quality comparison. Left to right: Our modified approach, Baseline Instant-NGP, Ground truth.

Static Assumption & Conditions: The architecture assumes static geometry and consistent lighting. Dynamic scenes or transparency violate the volume rendering assumptions. Reconstruction reliability degrades significantly unless scenes meet optimal conditions: moderate texture, diffuse materials, and 200-300 distinct viewpoints.

Computational Constraints: 50 credits of google cloud provided by the UMD was not enough. We purchased Google Colab Pro to train on our datasets. Although efficient, the training process still demands high-end GPUs. Resource constraints limited our ability to perform exhaustive hyperparameter tuning, suggesting potential for further optimization.

7. Conclusion

This work investigated the efficiency-quality trade-offs in neural rendering, validating Instant-NGP as the optimal baseline for practical deployment. By introducing adaptive hash grid resolution scheduling and an enhanced color MLP, we achieved a 0.66 dB PSNR improvement over the baseline while maintaining around 120-minute training time.

7.1. Key Findings

Efficiency vs. Quality: Our modifications prove that lightweight architectural changes can significantly reduce

artifacts without sacrificing speed. Adaptive scheduling successfully mitigates high-frequency noise and floaters, while the increased color network capacity captures view-dependent specular highlights that the baseline fails to resolve.

Geometric Consistency: Beyond photometric metrics, our approach yields superior 3D geometry and temporal stability in rendered videos, making it more suitable for downstream asset creation.

7.2. Future Directions

Generalization & Dynamics: Future work should explore pre-training on large-scale datasets (e.g., Objaverse) to enable few-shot reconstruction. Additionally, extending our adaptive scheduling to dynamic NeRF variants (e.g., HexPlane) could improve 4D reconstruction of moving objects.

Broader Impact: By achieving high-fidelity reconstruction on consumer-grade hardware in under 120 minutes, this work helps democratize 3D digitization, enabling applications in cultural preservation and e-commerce without requiring industrial compute clusters.

Acknowledgments

We thank instructors of this course Prof. George Zaki and Prof. Zachary Hanif for guidance and feedback throughout this project.

References

- [1] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, “NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis,” in *ECCV*, 2020.
- [2] T. Müller, A. Evans, C. Schied, and A. Keller, “Instant Neural Graphics Primitives with a Multiresolution Hash Encoding,” *ACM Trans. Graph.*, vol. 41, no. 4, 2022.
- [3] M. Tancik, E. Weber, E. Ng, R. Li, B. Yi, J. Kerr, T. Wang, A. Kristoffersen, J. Austin, K. Salahi, A. Ahuja, D. McAllister, and A. Kanazawa, “Nerfstudio: A Modular Framework for Neural Radiance Field Development,” in *SIGGRAPH Conference Proceedings*, 2023.
- [4] J. T. Barron, B. Mildenhall, M. Tancik, P. Hedman, R. Martin-Brualla, and P. P. Srinivasan, “Mip-NeRF: A Multiscale Representation for Anti-Aliasing Neural Radiance Fields,” in *ICCV*, 2021.
- [5] D. Verbin, P. Hedman, B. Mildenhall, T. Zickler, J. T. Barron, and P. P. Srinivasan, “Ref-NeRF: Structured View-Dependent Appearance for Neural Radiance Fields,” in *CVPR*, 2022.
- [6] J. L. Schönberger and J.-M. Frahm, “Structure-from-Motion Revisited,” in *CVPR*, 2016.