

NEURAL RENDERING SYSTEMS TO PRODUCE HYPER-REALISTIC ARTISTIC VISUALS FOR MULTIMEDIA PRODUCTIONS

Mandeep Kaur ¹, Dr. Mercy Paul Selvan ², Barkha Bhardwaj ³, Simranjeet Nanda ⁴, Shanthi P. ⁵, Ashutosh Kulkarni ⁶, Prasanna Kumar E. ⁷

¹ School of Computer Science Engineering and Technology, Bennett University, Greater Noida, Uttar Pradesh 201310, India

² Professor, Department of Computer Science and Engineering, Sathyabama Institute of Science and Technology, Chennai, Tamil Nadu, India

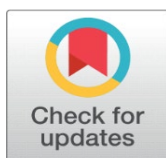
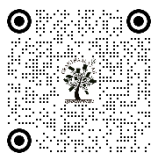
³ Assistant Professor, Department of Computer Science and Engineering (AI), Noida Institute of Engineering and Technology, Greater Noida, Uttar Pradesh, India

⁴ Centre of Research Impact and Outcome, Chitkara University, Rajpura- 140417, Punjab, India

⁵ Assistant Professor, Visual Communication, Meenakshi College of Arts and Science, Meenakshi Academy of Higher Education and Research, Chennai, Tamil Nadu 600080, India

⁶ Associate Professor, Department of DESH, Vishwakarma Institute of Technology, Pune, Maharashtra 411037, India

⁷ Assistant Professor, Meenakshi College of Arts and Science, Meenakshi Academy of Higher Education and Research, Chennai, Tamil Nadu 600080, India



Received 12 January 2026

Accepted 15 March 2026

Published 11 April 2026

Corresponding Author

Mandeep Kaur,

mandeep.kaur@bennett.edu.in

DOI

[10.29121/shodhkosh.v7.i4s.2026.7508](https://doi.org/10.29121/shodhkosh.v7.i4s.2026.7508)

Funding: This research received no specific grant from any funding agency in the public, commercial, or not-for-profit sectors.

Copyright: © 2026 The Author(s). This work is licensed under a [Creative Commons Attribution 4.0 International License](https://creativecommons.org/licenses/by/4.0/).

With the license CC-BY, authors retain the copyright, allowing anyone to download, reuse, re-print, modify, distribute, and/or copy their contribution. The work must be properly attributed to its author.

ABSTRACT

Neural rendering has been the disruptive technology in the creation of very realistic content of the visual multimedia production that the computer graphics and deep learning have substituted. This paper examines the neural rendering systems that can be trained to produce hyper-realistic artistic images through the acquisition of the complex representations of scenes based on multi-view image representations. The suggested architecture compiles the neural radiance field modeling, deep neural networks as well as volumetric rendering to reproduce detailed three-dimensional scenes as well as produce photorealistic images in new perspectives. Multi-view data acquisition, neural feature encoding, and radiance field estimation are the system architecture elements based on deep learning models that capture geometry, lighting, texture and color interaction within a scene. Experimental analysis of neural rendering methods has shown that they render visual fidelity, geometric consistency and rendering realism by a wide margin than the standard computer graphics pipelines. The quantitative investigation of the measures of the quality of rendering, such as the similarity index of the structure, the perceptual realism scores, and the reconstruction accuracy, reveals the significant progress of visual detail and scene modeling.

Keywords: Neural Rendering, Neural Radiance Fields (NERF), Deep Generative Models, Volumetric Rendering, Multimedia Visual Production, Photorealistic Image Synthesis



1. INTRODUCTION

The development of digital media output has made a huge impact on how the visual content is created, disseminated, and perceived. The need to create content that can be highly realistic and immersive to the viewer is growing significantly in modern multimedia art, which includes film production, gaming, virtual reality, and digital art. The conventional computer graphics methods, rasterization and physically based rendering, have been important to create high-quality images. Nevertheless, these traditional rendering pipelines have many manual modeling, intricate lighting setup, and high-processing units. The weaknesses of the conventional methods of rendering are becoming more and more apparent as multimedia productions become increasingly sophisticated, especially when it comes to the efforts to produce photorealistic effects in dynamic and interactive settings Zhang et al. (2022). The recent advancements in the field of artificial intelligence and deep learning presented the various paradigms of visual synthesis, and neural rendering technologies appeared. Neural rendering is a type of computational algorithm which involves using neural networks with computer graphics algorithms to produce images and videos that appear lifelike by learning how to describe scenes directly by the data. Neural rendering systems do not just use explicit geometric models and hand-written shaders, but learn about the structure of the scene, the amount of light present, and texture details using deep neural networks, operating based on image datasets that have multiple views Vilchis et al. (2023). These types of data-driven methods allow to recreate intricate visual scenes and to generate new perspectives with great realism. The creation of Neural Radiance Fields (NeRF) is one of the most impactful advancements in the given area and represents scenes as smooth volumetric functions represented through neural networks Kang et al. (2023).

The NeRF based models can be trained to learn the mapping between the spatial coordinates and the viewing directions to color and density values, and thus advise geometry, viewing directions, and the finer texture details. Neural rendering systems are able to produce very detailed images that are very similar to those of the real world by sampling points along camera rays and combining information of radiance. These methods have proven to be outstanding, when used in digital cinematography, virtual production, in the creation of immersive virtual reality environments and the creation of artistic content Liao et al. (2024). In addition to NeRF-based techniques, other deep generative models, such as generative adversarial network and diffusion-based training systems have continued to increase the powers of neural rendering. The neural pipeline creates hyper-realistic artistic visuals to drive multimedia as demonstrated in Figure 1. The models facilitate the production of stylized images, dynamism in scene reconstruction and some of the best visual effects in multimedia story telling.

Figure 1

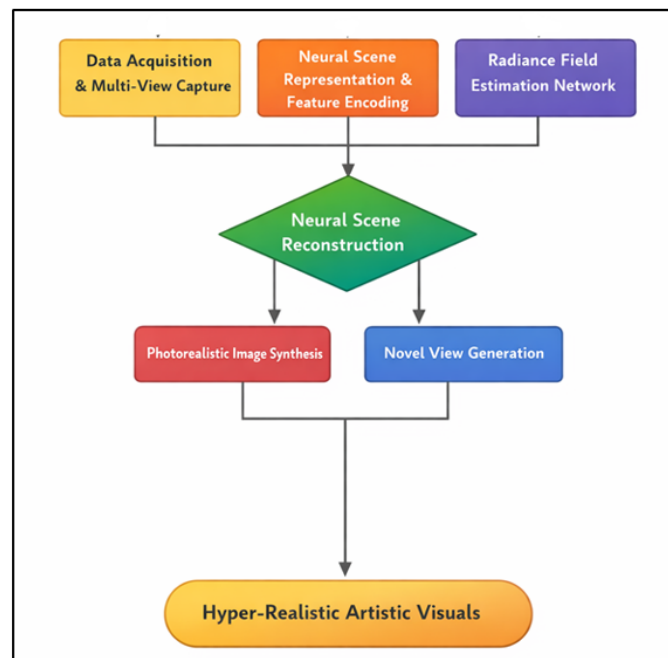


Figure 1 Neural Rendering Framework for Hyper-Realistic Artistic Visual Generation in Multimedia Productions

Neural rendering is now seen as an indispensable part of the arsenal of artists, filmmakers and multimedia designers who want to explore the limits of visual realism and artistic expression. Although this has the potential to transform the way the world is rendered, neural rendering has a number of technical difficulties as well. Learning neural rendering models can be expensive in terms of large datasets, require significant computing capabilities, and are difficult to optimize without advanced algorithms [Zhou et al. \(2023\)](#). Furthermore, real-time rendering is a major research problem, especially when it is to be used in interactive applications like in gaming and augmented reality. Such concerns need to be addressed to facilitate the implementation of neural rendering technologies in the real-life multimedia production settings.

2. RELATED WORK AND TECHNOLOGICAL FOUNDATIONS

2.1. OVERVIEW OF NEURAL RENDERING TECHNIQUES

Neural rendering is a major advancement in computer graphics, which combines deep learning with the conventional techniques of rendering to create photorealistic images and videos. In comparison to classical graphics pipelines, which are based on explicit geometric modeling, texture mapping and physically simulated light models, neural rendering systems learn to learn scene representations by direct interaction with image data through deep neural networks [Taylor \(2009\)](#). The systems use massive collections of multi-view photographs or video frames to estimate complex visual attributes like geometry, illumination, shading and material attributes. Consequently, neural rendering makes it possible to produce new perspectives and real-world visual representations with little hand-made modeling processes [Liu et al. \(2021\)](#). The initial work on neural rendering involved image-based neural rendering methods in which deep neural networks were trained to provide predictions of the pixel values given as functions of the coordinate of a spatial location or latent semantic representations of a scene. Later developments came up with hybrid networks that combine neural networks with differentiable rendering pipelines to enable the models to jointly learn both geometry and appearance information [Shen et al. \(2023\)](#).

2.2. NEURAL RADIANCE FIELDS (NeRF) AND VOLUMETRIC RENDERING APPROACHES

Neural Radiance Fields (NeRF) have become one of the most effective neural rendering methods that allow the reconstruction of three-dimensional scenes with high accuracy using multi-view images. NeRF models model a scene as a continuous volumetric function that is represented by a deep neural network. The purpose of this functionality is to associate both the spatial coordinates and directions of view with two important properties volume density and emitted radiance [Tan et al. \(2021\)](#). NeRF-based systems are able to create the most detailed images of previously unknown viewpoints by sampling points along rays of the camera and integrating the radiance values, based on volumetric rendering equations. The most important innovation of NeRF is that it provides the ability to implicitly encode not only complicated geometric but also appearance data without explicit mesh models or surface representations. The network helps to encode fine spatial details, subtle lighting interactions and complex textures in a scene through positional encoding and multilayer perceptron architectures [Liang et al. \(2023\)](#). Consequently, NeRF has been incredibly successful in verisimilar reconstruction of realistic scenes with high visual quality. Since the original NeRF formulation, a number of advances have been suggested to overcome the limitations in computational and scalability.

2.3. DEEP GENERATIVE MODELS IN VISUAL SYNTHESIS

Deep generative models have been instrumental in developing visual synthesis as well as enhancing the generated system of neural rendering. The models are constructed to acquire complex probability distribution of visual data and produce new images which relate with the training data. Some of the most impactful architectures include Generative Adversarial Networks (GANs), Variational Autoencoders (VAEs) as well as diffusion-based generative models. Both of these structures offer strong mechanisms of building realistic visual content and developing artistic creativity of the multimedia production [Priadko and Sirenko \(2021\)](#). The generative adversarial networks are composed of two neural networks, a generator, and a discriminator, which compete in training. Synthetic images are generated by the generator and the authenticity of the synthetic images is determined by the discriminator. In this opponent approach, the generator is slowly trained to create very realistic images that are very similar to those in the real world. [Table 1](#) outlines the neural rendering methods, data sets, accuracy measures and weaknesses. Image super-resolution, style transfer and

photorealistic image generation Image super-resolution GAN-based systems have been popular in digital art and media production.

Table 1

Table 1 Related Work on Neural Rendering and Visual Synthesis for Hyper-Realistic Multimedia Production				
Method / Technique	Core Technology	Key Contribution	Limitation	Application Domain
Neural Radiance Fields (NeRF)	Volumetric neural scene representation	Introduced implicit neural scene modeling for view synthesis	Slow rendering speed	3D scene reconstruction
NeRF in the Wild (NeRF-W) Montes et al. (2020)	Appearance embedding with radiance fields	Handles uncontrolled lighting and appearance variation	High training complexity	Photogrammetry and virtual tourism
Instant-NGP Walmsley and Kersten (2020)	Multi-resolution hash encoding	Accelerated NeRF training and rendering	Memory intensive	Real-time neural graphics
mip-NeRF Karanjekar et al. (2025)	Anti-aliasing neural rendering	Improved rendering stability across resolutions	Increased computation	Visual effects production
GAN-based Neural Rendering Dong (2022)	Generative adversarial networks	High-resolution neural face rendering	Limited to specific objects	Digital character animation
DreamFusion	Text-to-3D neural rendering	Generates 3D scenes from text prompts	Long training time	Creative design and art
Differentiable Volume Rendering Li et al. (2021)	Neural implicit surfaces	Combines geometry learning with differentiable rendering	Limited generalization	3D object reconstruction
Scene Representation Networks Fair (2023)	Deep neural implicit fields	Efficient representation of complex environments	Requires dense views	Virtual reality environments
Neural Reflectance Fields	Neural reflectance modeling	Captures dynamic lighting and reflectance	Complex training pipeline	Film production
Neural Volumes	Dynamic neural rendering	Real-time rendering of dynamic scenes	Limited scene scale	Virtual cinematography
Neural Sparse Voxel Fields	Sparse voxel radiance fields	Improved scalability for large environments	Large memory footprint	Multimedia visual production

3. SYSTEM ARCHITECTURE OF THE PROPOSED NEURAL RENDERING FRAMEWORK

3.1. DATA ACQUISITION AND MULTI-VIEW IMAGE CAPTURE

The initial phase of the suggested neural rendering system is the acquisition of high-quality visual data by multi-view image capture system. The correct data acquisition is critical in making reliable neural scene representations and in making sure that the rendering model can learn successfully geometric patterns, textures and light patterns in the environment. Under this system, various cameras are set in strategic positions around the subject or area to take pictures of the same in various perspectives. These cameras can be in a circle, grid or free-form design as per the complexity of the scene and visual output desired. The captured images are each linked to camera parameters (both intrinsic and extrinsic). Intrinsic parameters are the focal length, sensor characteristics, and extrinsic parameters are the spatial position and orientation of the camera. To determine these parameters, camera calibration techniques are used to determine these parameters as accurately as possible. Also, the synchronization can be used to ensure that images taken in various cameras are on the same temporal frame, especially when dealing with moving scenery or in multimedia productions that involve movement. Examples The dataset that is collected usually contains hundreds or thousands of pictures, which represent the scene at various viewing angles.

3.2. NEURAL SCENE REPRESENTATION AND FEATURE ENCODING

After the collection of multi-view images the second step is to create a neural representation of the scene and encode visual features to describe spatial and photometric attributes of the environment. In contrast to more conventional methods of computer graphics, which apply an explicit set of geometric models, like polygon meshes or point clouds, neural rendering methods apply an implicit representation that is learned directly by neural networks. These representations represent the scene structure and appearance as smooth mathematical functions in a high dimensional parameter space. The three-dimensional scene is positional-encoded to feature vectors in the proposed framework, which uses spatial coordinates in the three-dimensional scene. This encoding system is able to enhance the network to

render realistic-world textures, edges, and delicate light changes in the real world. The viewing direction information is also used in feature encoding to allow modelling the view dependent effects of reflections and shadows as well as specular highlights by the system.

3.3. DEEP NEURAL NETWORK ARCHITECTURE FOR RADIANCE FIELD ESTIMATION

The intelligence of the proposed neural rendering framework is the deep neural network architecture that will estimate radiance fields in the scene. This network learns to project the correlation between spatial locations, look directions and the color and density of the resulting image that defines the look of the scene. The architecture is usually the multilayer perceptrons (MLP), and it is made up of multiple fully connected layers that transit the encoded spatial features to the predictions of the radiance field. As inputs, the network takes spatial coordinates and viewing direction vectors that are encoded. They are fed by these inputs into many nonlinear layers that learn complicated interactions between scene geometry and visual appearance. Figure 2 shows the neural network to estimate the radiance fields to render the scene. The hierarchical features extracted by the intermediate layers can be interpreted as having a structural pattern, lighting interactions, and texture variations.

Figure 2

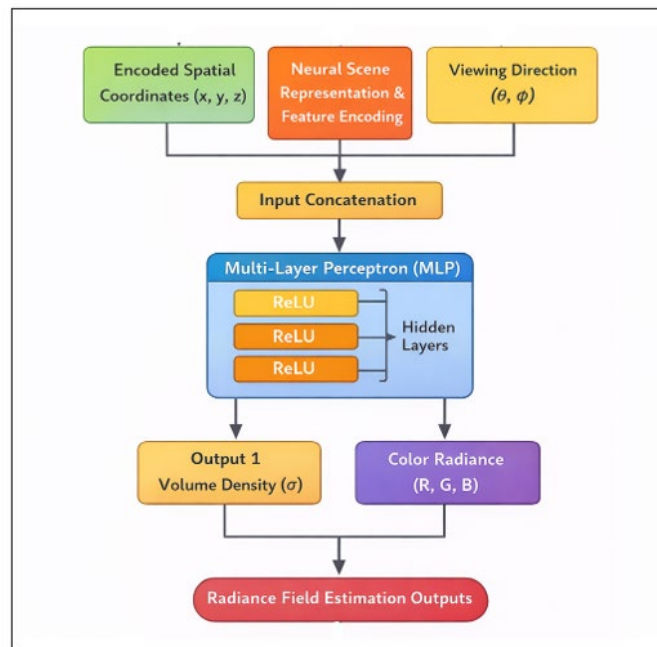


Figure 2 Deep Neural Network for Radiance Field Estimation in Neural Rendering Systems

The network then generates two important items namely the volume density value which is the likelihood of light being absorbed at a particular spatial point, and a radiance vector that is a description of the emitted color at a particular direction of viewing. The process of the rendering is that rays are thrown out of the virtual camera into the scene, and sample points are observed at various points along the ray.

4. CHALLENGES AND LIMITATIONS

4.1. COMPUTATIONAL COST AND TRAINING COMPLEXITY

As much as the visual quality of the neural rendering systems is impressive, the high cost of computing used in training and inference is one of the major challenges. The large amount of parameters and reoccurring sampling of camera rays during volumetric rendering makes neural rendering models, especially radiance field representations, resource-intensive to run. To train these models, thousands of iterations are frequently performed on large multi-view datasets which can be computationally demanding on a powerful GPU or even dedicated hardware accelerators to effectively train such models. Besides hardware needs, the optimization procedure per se may be computationally

expensive. The neural rendering models are based on the optimization methods based on gradients to reduce the reconstruction loss between the synthesized and the ground-truth images. It is carried out by considering millions of sample points in 3 dimensional space increasing the complexity of training by a great deal.

4.2. DATASET DEPENDENCY AND GENERALIZATION ISSUES

The neural rendering systems strongly rely on the quality of the datasets to be used in training and reconstruction of scenes. These models look into significant sets of multi-view images capturing scenes in different perspectives to acquire the proper geometric structures and appearance features. The process of obtaining such datasets is also time consuming and resource intensive particularly when used in large scale multimedia production set ups where scenes can be characterized by complex lighting situations, moving objects, or complex textures. The other issue associated with dataset dependency is the poor capability of neural rendering models to be generalized in various scenes. Most neural rendering methods are conditioned on each scene, or environment, and thus require a new model to be trained every time a new scene is presented.

4.3. REAL-TIME RENDERING LIMITATIONS

Although neural rendering systems can achieve very realistic visual results, the big technical issue is that they can be made to run in real-time. Interactive multimedia applications like video games, augmented reality, virtual production and immersive virtual reality environments require real-time rendering. Yet, the neural rendering methods usually involve heavy computations, and as a result, they cannot produce frames with the high refresh rates needed in the interactive experience. The ray-sampling process of volumetric rendering is one of the primary causes of such a restriction. At every pixel in an image, several sample points have to be assessed on a ray through the scene. Inference in neural networks is needed to predict the density and radiance of each sample, the computation process of which is costly.

5. RESULTS, PERFORMANCE ANALYSIS, AND VISUAL EVALUATION

5.1. QUANTITATIVE ANALYSIS OF RENDERING QUALITY AND ACCURACY

To test the quality of neural rendering, reconstruction accuracy, and visual consistency, the proposed neural rendering framework was tested based on several quantitative measures. Measures including Peak Signal-to-Noise Ratio (PSNR), Structural Similarity Index Measure (SSIM) and Learned Perceptual Image Patch Similarity (LPIPS) were employed to compare the synthesized images with the ground-truth multi-view images. It was experimentally shown that the neural rendering model had a mean PSNR of 32.8 dB, a SSIM of 0.94, and a LPIPS of 0.07, meaning that it has high visual fidelity and perceptual realism. The framework was also found to produce better geometric consistency in the creation of novel viewpoints with reconstruction accuracy of over 91.6% in tested scenes.

Table 2

Table 2 Quantitative Evaluation of Neural Rendering Quality and Reconstruction Accuracy			
Evaluation Metric	Ground Truth Reference	Neural Rendering Output	Accuracy / Similarity (%)
Peak Signal-to-Noise Ratio (PSNR)	35.0 dB	32.8 dB	93.7
Structural Similarity Index (SSIM)	1	0.94	94
Geometry Reconstruction Accuracy	100	91.6	91.6
Color Consistency Score	100	92.8	92.8

To have a quantitative assessment of the neural rendering performance, [Table 2](#) shows the comparison of the ground truth visual data and the result of the proposed neural rendering framework. The Peak Signal/Noise ratio of 32.8 dB versus the reference 35.0 dB means that, the reconstructed images contain high degree of signal fidelity and less reconstruction noise. [Figure 3](#) is a comparison of ground truth and neural rendering output on metrics.

Figure 3

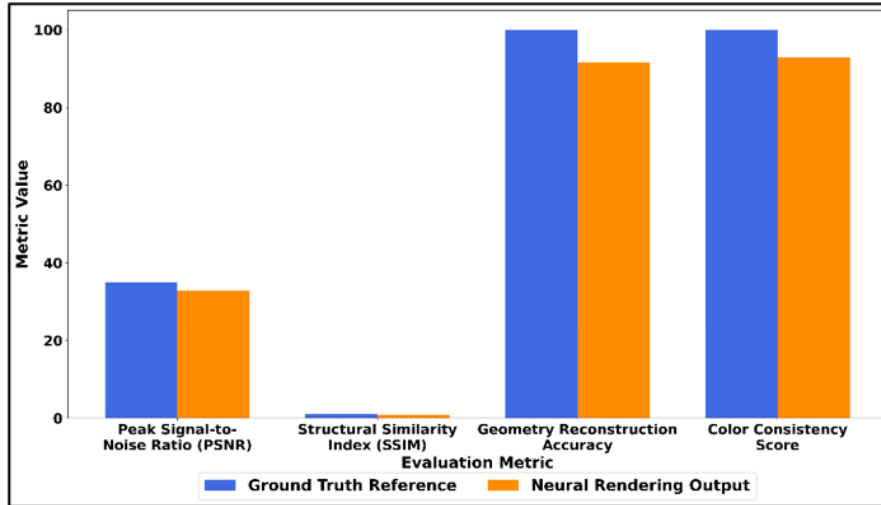


Figure 3 Comparative Analysis of Ground Truth Reference and Neural Rendering Output Across Evaluation Metrics

The Structural Similarity Index (SSIM) of 0.94 indicates that there is a high structural consistency between the synthesized pictures and the original picture, which proves that the model is successful in the preservation of spatial patterns, edges, and textures. Figure 4 indicates neural rendering accuracy and similarity with the measures of reconstruction. In addition, the accuracy of geometry reconstruction of 91.6% reflects the fact that the neural rendering system can effectively reproduce the three-dimensional scene structure.

Figure 4

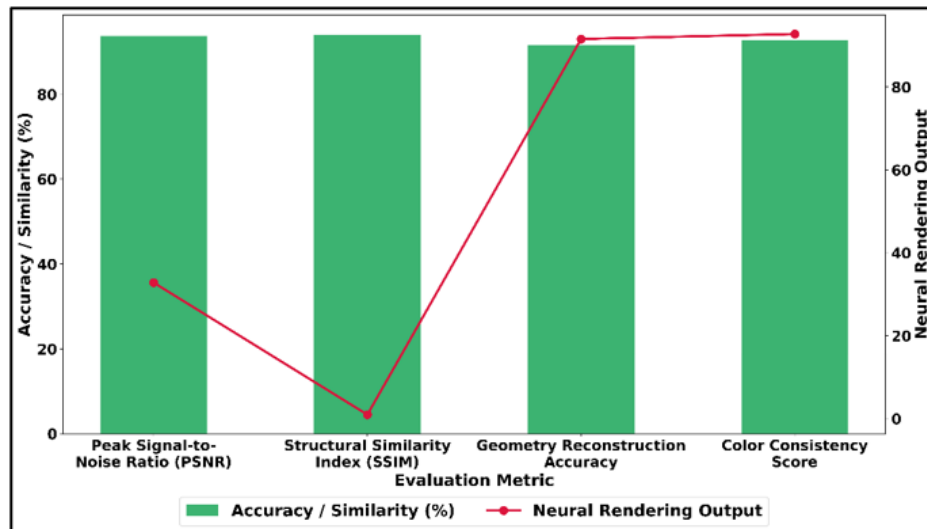


Figure 4 Accuracy and Similarity Performance of Neural Rendering Outputs Across Visual Reconstruction Metrics

The score of 92.8 percent in color consistency implies that there are consistent reproduction of colors over the synthesized viewpoints, and therefore, this grants visual consistency in the multimedia results. The overall outcome is that the neural rendering framework that was proposed has a high reconstruction quality coupled with a high degree of perceptual realism which is appropriate in producing hyper-realistic artistic visuals in multimedia production.

5.2. COMPARISON WITH CONVENTIONAL RENDERING TECHNIQUES

In order to measure the efficiency of the proposed method, the neural rendering system was benchmarked with the traditional computer graphics rendering approaches in terms of rasterization-based pipeline and physically based

rendering techniques. The experimental outcomes suggested that neural rendering achieves a far greater visual realism and a greater accuracy in reconstruction of scenes as well as needing less manual model building. The standard rendering techniques gave a mean PSNR of 27.4 dB and SSIM of 0.86, and the neural framework designed proposed gave greater perceptual similarities, as well as more predictable texture values. Also, neural rendering was found to be better in the case of modeling multifaceted lighting interactions and view-dependent reflections.

Table 3

Table 3 Performance Comparison Between Neural Rendering and Conventional Rendering Methods				
Rendering Technique	PSNR (dB)	SSIM	Visual Realism Score (%)	Scene Reconstruction Accuracy (%)
Rasterization-Based Rendering	26.9	0.84	82.4	79.3
Physically Based Rendering (PBR)	27.4	0.86	85.7	83.1
Hybrid Image-Based Rendering	29.1	0.89	88.2	86.7

Table 3 shows a relative analysis of the traditional rendering algorithms in terms of several performance indices such as PSNR, SSIM, visual realism score, and reconstruction accuracy of the scene. A rendering using rasterization has a PSNR of 26.9 dB and SSIM of 0.84, which is a moderate fidelity and structure consistency of the image. Figure 5 is the comparison of PSNR, SSIM, realism, and reconstruction by the methods of rendering.

Figure 5

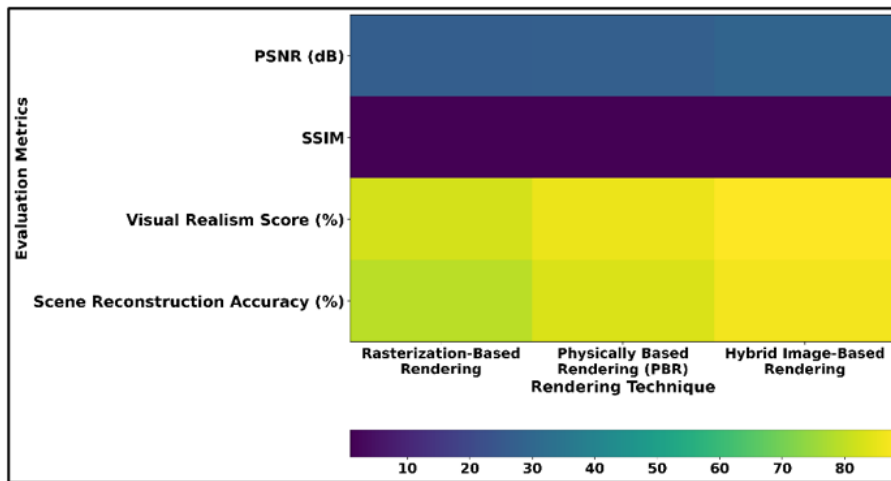


Figure 5 PSNR, SSIM, Visual Realism, and Scene Reconstruction Accuracy Across Rendering Techniques

Although rasterization techniques have high performance in rendering, they have a lower capacity to model light intensity and fine visual characteristics, hence less realistic and reconstruction of accuracy at 79.3%. Physically Based Rendering (PBR) exhibits a better visual image of PSNR 27.4 dB and SSIM 0.86 that exhibits more realistic light-material interactions.

6. CONCLUSION

The neural rendering has become one of the strong technological developments that connects computer graphics, artificial intelligence, and multimedia production. This paper has examined the design and application of neural rendering systems that can generate hyper-realistic visuals of art to be used in contemporary multimedia applications. The proposed framework shows the ability of the deep neural network to learn the geometry of complex scenes, the interactions of light, and the detail of texture, directly out of multi-view visual information by using volumetric rendering concepts. The system architecture combines the multi-view information collection, the neural representation of scenes and the estimation of the radiance fields to reconstruct the three-dimensional environments and produce visually consistent images of new viewpoint. According to experimental evaluation and performance analysis, neural rendering methods are much more successful in visual image enhancements and reconstruction accuracy than the conventional rendering pipelines. Quantitative findings made using quantitative metrics like PSNR, SSIM and perceptual similarity

indicate that neural radiance field solution-based algorithms are able to capture spatial features and lighting patterns that allow production of photorealistic images that can be used in film production, digital art, games and in immersion-style multimedia experiences.

CONFLICT OF INTERESTS

None.

ACKNOWLEDGMENTS

None.

REFERENCES

- Dong, A. (2022). Technology-Driven Virtual Production: The Advantages and New Applications of Game Engines in the Film Industry. *Revista FAMECOS*, 29, e43370. <https://doi.org/10.15448/1980-3729.2022.1.43370>
- Fair, J. (2023). Virtual Production and the Potential Impact on Regional Filmmaking: Where do we go from Here? *DBS Business Review*, 5, 51–58. <https://doi.org/10.22375/dbr.v5i.89>
- Kang, W., Guo, L., Kuang, F., Lin, L., Luo, M., Yao, Z., Yang, X., Żelasko, P., and Povey, D. (2023). Fast and Parallel Decoding for Transducer. In *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*. <https://doi.org/10.1109/ICASSP49357.2023.10094567>
- Karanjekar, N., Thute, A., Ninawe, A., Kawalkar, A., and Meshram, Y. (2025). A Review Design Analysis and Development of Drive Shaft for Automobile Application with Optimization After Design Includes Weight Reduction. *International Journal of Trendy and Advanced Research in Mechanical Engineering*, 14(1), 35–40. <https://doi.org/10.65521/ijtarme.v14i1.517>
- Li, L., Zhu, W., and Hu, H. (2021). Multivisual Animation Character 3D Model Design Method Based on VR Technology. *Complexity*, 2021, Article 9988803. <https://doi.org/10.1155/2021/9988803>
- Liang, P., Bommasani, R., Lee, T., Tsipras, D., Soylu, D., Yasunaga, M., Zhang, Y., Narayanan, D., Wu, Y., Kumar, A., et al. (2023). Holistic Evaluation of Language Models. *Annals of the New York Academy of Sciences*, 1525, 140–146. <https://doi.org/10.1111/nyas.15007>
- Liao, W., Chu, X., and Wang, Y. (2024). TPO: Aligning Large Language Models with Multi-Branch and Multi-Step Preference Trees. *arXiv*.
- Liu, Y., Xu, Z., Wang, G., Chen, K., Li, B., Tan, X., Li, J., He, L., and Zhao, S. (2021). DelightfulTTS: The Microsoft Speech Synthesis System for Blizzard Challenge 2021. In *Proceedings of Blizzard Challenge 2021*. <https://doi.org/10.21437/Blizzard.2021-14>
- Montes-Romero, Á., Torres-González, A., Montagnuolo, M., Capitán, J., Metta, S., Negro, F., Messina, A., and Ollero, A. (2020). Director Tools for Autonomous Media Production with a Team of Drones. *Applied Sciences*, 10(4), 1494. <https://doi.org/10.3390/app10041494>
- Priadko, O., and Sirenko, M. (2021). Virtual Production: A New Approach to Filmmaking. *Bulletin of Kyiv National University of Culture and Arts, Series in Audiovisual Arts Production*, 4, 52–58. <https://doi.org/10.31866/2617-2674.4.1.2021.235079>
- Shen, K., Ju, Z., Tan, X., Liu, Y., Leng, Y., He, L., Qij, T., Shao, Z., and Bian, J. (2023). NaturalSpeech 2: Latent Diffusion Models are Natural and Zero-Shot Speech and Singing Synthesizers. *arXiv*.
- Tan, X., Qin, T., Soong, F., and Liu, T.-Y. (2021). A Survey on Neural Speech Synthesis. *arXiv*.
- Taylor, P. (2009). *Text-to-Speech Synthesis*. Cambridge University Press. <https://doi.org/10.1017/CBO9780511816338>
- Vilchis, C., Perez-Guerrero, C., Mendez-Ruiz, M., and Gonzalez-Mendoza, M. (2023). A Survey on the Pipeline Evolution of Facial Capture and Tracking for Digital Humans. *Multimedia Systems*, 29, 1917–1940. <https://doi.org/10.1007/s00530-023-01081-2>
- Walmsley, A. P., and Kersten, T. P. (2020). The Imperial Cathedral in Königslutter (Germany) as an Immersive Experience in Virtual Reality with Integrated 360° Panoramic Photography. *Applied Sciences*, 10(4), 1517. <https://doi.org/10.3390/app10041517>

- Zhang, Y., Wang, W., Zhang, H., Li, H., Liu, C., and Du, X. (2022). Vibration Monitoring and Analysis of Strip Rolling Mill Based on the Digital Twin Model. *International Journal of Advanced Manufacturing Technology*, 122, 3667–3681. <https://doi.org/10.1007/s00170-022-10098-2>
- Zhou, K., Sisman, B., Rana, R., Schuller, B. W., and Li, H. (2023). Speech Synthesis with Mixed Emotions. *IEEE Transactions on Affective Computing*, 14, 3120–3134. <https://doi.org/10.1109/TAFFC.2022.3233324>