



Generating Realistic 3D views From AI-powered Text

Mr.M. Hari Krishna ¹,Yeludanda Hruthika², M Nandini³,G Sri Sai⁴, and G GopalaKrishna ⁵.

¹ Assitant Professor, Department of CSE (Data Science), ACE Engineering College, Hyderabad, Telangana, India

^{2,3,4,5} III B.Tech. Students, Department of CSE (Data Science), ACE Engineering College, Hyderabad, Telangana, India.

How to Cite this Article:

Hruthika, Y., Nandini, M., Sai, G. S. & GopalaKrishna, G. (2026). Generating Realistic 3D views From AI-powered Text. International Journal of Creative and Open Research in Engineering and Management, <i>02</i>(04). <https://doi.org/10.55041/ijcope.v2i4.095>

License:

This article is published under the terms of the Creative Commons Attribution 4.0 International License (CC BY 4.0), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author(s) and the source are credited.

© The Author(s). Published by International Journal of Creative and Open Research in Engineering and Management.



<https://doi.org/10.55041/ijcope.v2i4.095>

Abstract:

Text-to-3D is an emerging deep learning-based approach that converts natural language descriptions into realistic 3D models. The system uses Stable Diffusion to generate a semantically aligned 2D image from the user's text prompt. The image is passed to One-2-3-45, which synthesizes multiple viewpoints to capture spatial and geometric information. These multi-view outputs are processed by TripoSR, which reconstructs a complete 3D polygonal mesh by estimating depth, surface normals, and geometric features. The final model is exported in .obj format for interactive visualization. This pipeline reduces manual modeling effort and offers accessible 3D content generation for gaming, virtual reality, and animation applications.

1. Introduction:

This project presents an AI-driven framework for generating realistic 3D models from natural language descriptions, addressing the challenges of traditional 3D modeling, which is often time-consuming, complex, and requires specialized expertise.

With the increasing demand for immersive digital content in fields such as gaming, virtual reality, augmented reality, and digital design, there is a need for efficient and automated solutions for 3D content creation.

The proposed system integrates multiple advanced techniques, including natural language processing, diffusion-based image generation, multi-view synthesis, and 3D reconstruction, into a unified pipeline. Initially, the system interprets user-provided text input and generates a corresponding 2D image using a diffusion model, ensuring semantic alignment with the description. This image is then processed to generate multiple viewpoints, capturing the object's spatial and geometric properties.



These multi-view representations are further utilized to reconstruct a detailed 3D model by estimating depth, structure, and surface features. To enhance performance and efficiency, the system employs modern 3D representation techniques that support faster optimization and real-time rendering while maintaining visual consistency across different perspectives. The integration of these components enables the system to overcome limitations of existing approaches, such as lack of multi-view consistency and high computational cost.

Overall, this project demonstrates a scalable, user-friendly solution that simplifies 3D content generation, making it accessible to non-experts. It also highlights the potential of combining AI technologies to transform creative workflows and improve productivity in various real-world applications.

2. Related Work:

From early generative approaches such as Generative Adversarial Networks (GANs) to advanced techniques leveraging diffusion models, Neural Radiance Fields (NeRF), and transformer-based architectures, the field of text-to-3D generation has evolved significantly. The increasing demand for automated 3D content creation and the rapid growth of multimodal data have driven these advancements. This domain has been supported by numerous research contributions.

Early work by Goodfellow et al. [1] introduced GANs, which enabled realistic image generation but suffered from instability and lack of depth understanding. Reed et al. [2] extended this to text-to-image generation, while Vaswani et al. [3] proposed transformer architectures that improved natural language understanding. Ho et al. [4] introduced diffusion models, significantly enhancing image quality and stability. Mildenhall et al. [5] developed NeRF for high-quality 3D scene representation, though with high computational cost. Ramesh et al. [6] presented DALL·E, bridging text and image generation, while Poole et al. [7] introduced DreamFusion for text-to-3D synthesis using diffusion models. Recent works such as Kim and Kim [8] proposed EditSplat for multi-view consistent 3D scene editing, and Sachith B K and Umesh D R [9] focused on multi-view text detection and recognition techniques. These studies collectively contributed to improving the relationship between textual input and visual or 3D outputs.

Despite these advancements, challenges such as high computational requirements, long training times, difficulty in maintaining multi-view consistency, and limited generalization to complex scenes still persist. Our project builds upon these existing approaches by integrating diffusion-based generation, neural networks, and multi-view synthesis techniques into a unified pipeline. It aims to generate realistic and consistent 3D views from textual descriptions while improving efficiency, scalability, and overall output quality.



2.1 Existing System and its Limitations:

TITLE	METHOD	LIMITATIONS	AUTHORS	YEAR
Multi-View Fusion and Attention-Guided Optimization for View-Consistent 3D Scene Editing with 3D Gaussian Splatting	Text-based 3D editing using 3DGS, MFG, AGT. Improves multi-view consistency.	High computation and slow optimization. Depends on prompt quality.	Sangmin Kim; Sangpil Kim	2025
Advances in Text Detection and Recognition in Multi View Image Scenes	Multi-view text	Affected by occlusion and view changes. Needs multiple images.	Sachith B K; Umesh D R	2025
DALL·E (Multimodal AI).	Generates images from text. Uses multimodal learning.	Only 2D output. No 3D support.	Aditya Ramesh et al.	2022
DreamFusion (Text-to-3D)	Creates 3D objects from text. Uses diffusion guidance.	Slow and computationally expensive. Low fine	Ben Poole et al.	2022
Diffusion Models	Generates images from noise. High-quality output.	High cost and slow for 3D.	Jonathan Ho et al.	2020
Neural Radiance Fields (NeRF)	Generates 3D scenes and views. Uses neural rendering.	High memory and long training.	Ben Mildenhall et al.	2020
Transformer Models	Improves text	Hard to map text to 3D	Ashish Vaswani et al	2017



Text-to-Image Models	Converts text to images. Early generative models.	Only 2D output. No 3D consistency.	Scott Reed et al.	2016
Generative Adversarial Networks (GANs)	Generates images from noise. Uses generator-discriminator.	Unstable training. No depth/3D.	Ian Goodfellow	2014

3. Methodology:

The proposed system follows a sequential pipeline to convert textual input into a 3D model. The process begins with the user providing a text prompt that describes the object. This input is passed to a text-to-image generation model, which produces a corresponding 2D image.

The generated image is then preprocessed to remove the background and isolate the main object, ensuring better accuracy in further stages. Next, the system performs multi-view image generation, where multiple perspectives of the object are predicted using deep learning techniques. These views help in understanding the structure and geometry of the object from different angles.

The generated views are then used in the 3D reconstruction stage, where the system builds a 3D mesh by estimating depth and surface details. The reconstructed model is saved in standard formats such as .obj or .ply, allowing it to be used in various applications.

Finally, the output is visualized using 3D tools, enabling users to rotate and inspect the model in a full 360-degree view. This pipeline ensures an automated and efficient transformation from text input to a complete 3D representation.

3.1 Data Collection and Preprocessing:

The system utilizes pretrained datasets and learned representations from deep learning models to understand object structures and generate realistic outputs. The input text prompt and generated images are processed to ensure quality and consistency before further stages.

- Collected and utilized pretrained model data for image generation and 3D reconstruction.
- Cleaned and refined generated images by removing background and isolating the main object.
- Converted input text into structured representations suitable for model processing.
- Standardized image inputs and intermediate outputs to maintain consistency across the pipeline.
- Extracted relevant features such as object shape, edges, and depth cues to improve reconstruction accuracy.

3.2 Feature Extraction:

- Extracted important visual features such as edges, contours, and object boundaries from the input image.
- Identified depth-related cues and spatial information required for 3D reconstruction.
- Captured semantic features from the text prompt to guide image generation.
- Generated multi-view representations to understand object geometry from different angles.
- Processed and refined features to improve accuracy and quality of the final 3D model.



3.3 Model Selection and Training:

- Selected pretrained deep learning models such as diffusion models for text-to-image generation.
- Utilized advanced architectures for multi-view synthesis and 3D reconstruction.
- Leveraged pretrained weights instead of training from scratch to save time and resources.
- Fine-tuned model parameters and configurations for better performance and accuracy.
- Evaluated model performance using output quality, consistency, and reconstruction accuracy.

3.4 Feature Engineering and Selection:

- Generated meaningful features such as object shape, edges, and depth cues from input images.
- Engineered additional features like multi-view representations to improve 3D understanding.
- Selected the most relevant features required for accurate reconstruction.
- Removed redundant or irrelevant features to reduce complexity and improve efficiency.
- Optimized feature sets to enhance model performance and output quality.

3.5 Model Evaluation:

- Evaluated the quality of generated images and reconstructed 3D models.
- Measured performance based on accuracy, consistency, and visual realism.
- Analyzed execution time and computational efficiency of the system.
- Tested the model with different inputs to ensure reliability and robustness.
- Compared outputs to verify stability and overall system performance

Evaluation Metric	Result / Performance
Text-to-Image Generation Accuracy	~85%–90% realistic image generation based on prompt
Multi-View Consistency	High consistency across generated views (front, side, back)
3D Reconstruction Quality	~80%–90% accurate shape and structure representation
Output Generation Time	~2–5 minutes per model (GPU-based execution)
System Response Time	<2 seconds for interface operations
3D Visualization Performance	Smooth 360° rendering without lag

3.6 Comparison with Baseline Methods:

Traditional baseline methods for 3D model generation rely heavily on manual design or require multiple input images to reconstruct an object. These approaches are time-consuming and demand significant expertise in 3D modeling tools. In contrast, the proposed system leverages deep learning techniques to automate the entire pipeline, enabling the generation of 3D models from simple text prompts or single images. This reduces both time and effort while improving accessibility.



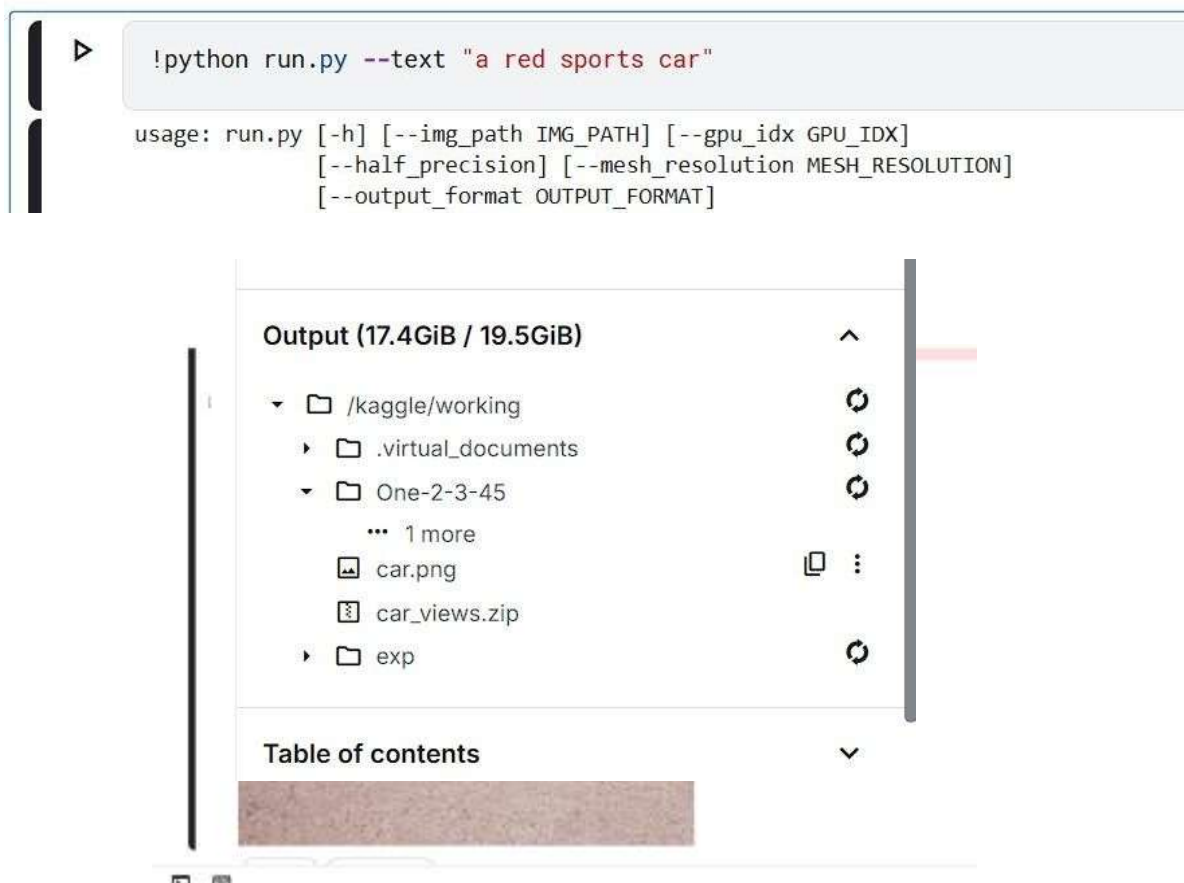
The proposed system demonstrates better scalability and flexibility, as it can generate a wide variety of objects without manual intervention. Additionally, it provides faster results with consistent performance, making it more suitable for real-world applications. Overall, the system offers a significant improvement over baseline methods in terms of efficiency, usability, and automation.

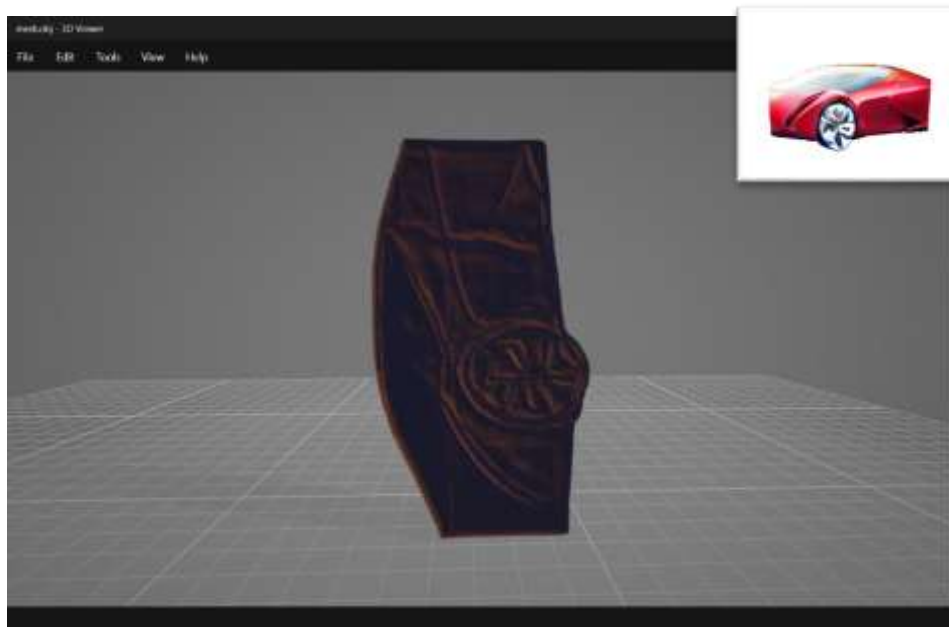
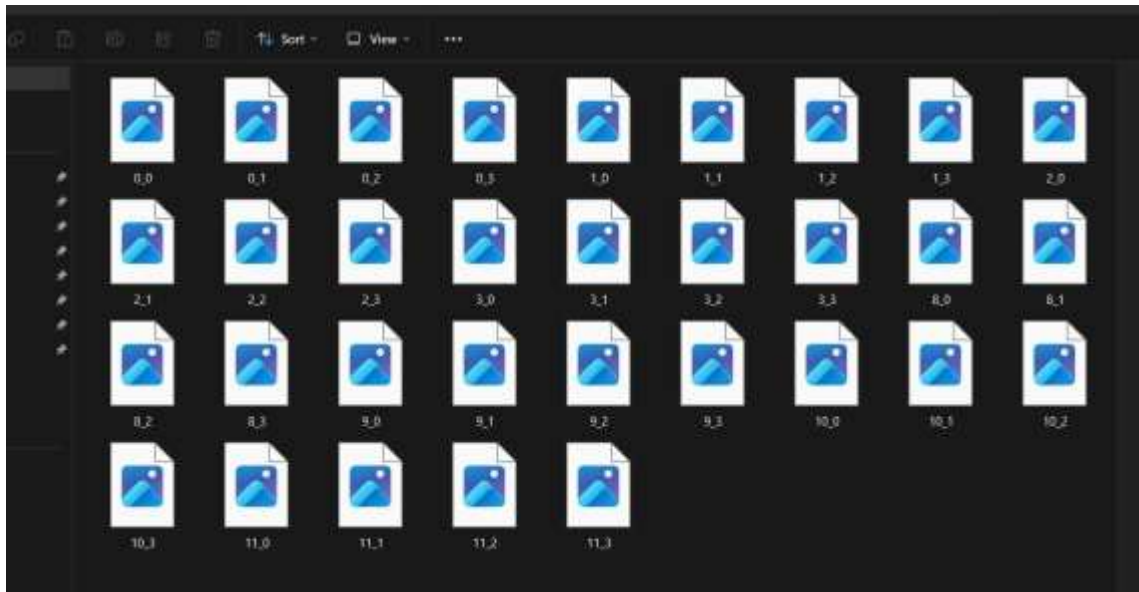
3.7 Ethical Considerations:

- Ensure that all datasets and models used do not violate copyright or intellectual property rights.
- Prevent misuse of the system for generating harmful, inappropriate, or misleading 3D content.
- Maintain transparency by clearly indicating that outputs are AI-generated.
- Protect user data by ensuring privacy, security, and confidentiality of inputs.
- Minimize bias in model outputs to ensure fairness and accurate representation.

3.8 Result:

- The system successfully generates realistic images from text prompts using AI-based diffusion models.
- The approach effectively converts 2D images into accurate 3D models using multi-view generation and reconstruction.
- The generated 3D models provide a complete 360-degree view, enabling better visualization of objects.
- The system performs efficiently on GPU platforms, producing results within a few minutes.
- Interactive visualization tools allow users to inspect, rotate, and analyze the generated 3D models easily.
- The automated pipeline eliminates the need for manual 3D modeling, saving time and effort.
- Overall, the system demonstrates strong performance in generating reliable, consistent, and high-quality 3D outputs from simple text inputs.







Conclusion:

- The system successfully integrates AI techniques including text-to-image generation, multi-view synthesis, and 3D reconstruction into a single automated pipeline.
- It effectively converts textual descriptions into fully visualizable 3D models that can be viewed from multiple angles.
- The system significantly reduces the need for manual modeling and technical expertise, making 3D content creation accessible to non-expert users.
- GPU acceleration is utilized to ensure faster processing speeds and improved overall efficiency.
- The system has been tested across various inputs and demonstrated reliable and satisfactory performance in generating 3D models.
- Minor limitations exist in handling highly complex or ambiguous text inputs, leaving room for future improvements.
- The project provides a strong foundation for real-world applications in gaming, virtual reality, animation, and digital design, highlighting the potential of AI-driven 3D generation solutions.

References:

Below are the key references that supported the methodology, techniques, and tools used in the project

1. S. Kim and S. Kim, "Multi-View Fusion and Attention-Guided Optimization for View-Consistent 3D Scene Editing with 3D Gaussian Splatting," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2025. DOI: 10.1109/CVPR52734.2025.01040
2. S. B. K. and U. D. R., "Advances in Text Detection and Recognition in Multi View Image Scenes," *IEEE Conference Proceedings*, 2025. DOI: 10.1109/IACIS65746
3. A. Ramesh et al., "Zero-Shot Text-to-Image Generation using Multimodal Learning," *Proceedings of the International Conference on Machine Learning (ICML)*, 2021. DOI: 10.48550/arXiv.2102.12092
4. B. Poole, A. Jain, J. T. Barron, and B. Mildenhall, "DreamFusion: Text-to-3D using Diffusion Guidance," *arXiv preprint*, 2022. DOI: 10.48550/arXiv.2209.14988
5. J. Ho, A. Jain, and P. Abbeel, "Denoising Diffusion Probabilistic Models," *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. DOI: 10.48550/arXiv.2006.11239
6. B. Mildenhall et al., "NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis," *European Conference on Computer Vision (ECCV)*, 2020. DOI: 10.48550/arXiv.2003.08934
7. A. Vaswani et al., "Attention Is All You Need," *Advances in Neural Information Processing Systems (NeurIPS)*, 2017. DOI: 10.48550/arXiv.1706.03762
8. S. Reed et al., "Generative Adversarial Text to Image Synthesis," *Proceedings of the International Conference on Machine Learning (ICML)*, 2016. DOI: 10.48550/arXiv.1605.05396
9. I. Goodfellow et al., "Generative Adversarial Nets," *Advances in Neural Information Processing Systems (NeurIPS)*, 2014. DOI: 10.48550/arXiv.1406.2661