

International Journal of Advanced Research in Science, Communication and Technology

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Volume 5, Issue 1, July 2025



VisionGen 3D: A Deep Learning Framework for Animation and Image Enhancement

Sanskruti P. Upadhyay¹, Noor Siddiqui², Pranay Vitekar³, Mahek Pathan⁴, Pushpa Tandekar⁵

Student, Computer Science & Engineering, Shri Sai College of Engineering & Technology, Bhadrawati, India¹²³⁴ Asst. Professor, Computer Science & Engineering, Shri Sai College of Engineering & Technology, Bhadrawati, India⁵

Abstract: The field of 3D animation is experiencing a significant transformation driven by deep learning advancements. This paper introduces a robust and scalable system that automates 3D facial animation from a single static image, utilizing Generative Adversarial Networks (GANs) and Variational Autoencoders (VAEs). The system is optimized for real-time deployment on cloud platforms, eliminating the need for manual animation techniques or motion-capture equipment. It accepts a static facial image along with a driving video to accurately map expressions and movements with high visual fidelity. Quantitative evaluation demonstrates strong performance, with a Structural Similarity Index (SSIM) of 0.87, Fréchet Inception Distance (FID) of 23.4, and real-time processing at 23 frames per second (FPS). The proposed framework supports a wide range of facial types—including human, cartoon, and avatar faces—offering a generalized and accessible solution for 3D animation generation through deep learning.

Keywords: 3D Animation, Deep Learning, GAN, VAE, Real-Time Processing, SSIM, Facial Animation

I. INTRODUCTION

In the current era of artificial intelligence, traditional animation workflows are undergoing a profound transformation. Conventional manual animation methods are not only time-intensive but also demand skilled expertise and substantial resources. To address these limitations, this paper presents an AI-driven framework that leverages deep learning models—specifically Generative Adversarial Networks (GANs) and Variational Autoencoders (VAEs)—to generate 3D facial animations directly from static images. The primary objective is to democratize the production of high-quality animations through a real-time, cloud-based system that is accessible even to users without specialized technical knowledge.

II. LITERATURE REVIEW

Research in AI-driven animation has evolved significantly, transitioning from early morphing-based methods to sophisticated, data-driven models capable of producing highly realistic renderings. Early work such as the **First-Order Motion Model** (Siarohin et al., 2019) introduced motion transfer through unsupervised keypoint detection, laying the foundation for image animation from a single frame. Subsequent models like **Liquid Warping GAN** improved generalization across variations in identity and pose. **Few-shot adversarial models** (Zakharov et al., 2019) demonstrated effective results with limited training data, enabling rapid adaptation to new identities. Despite these advancements, challenges such as temporal consistency and full-body animation remain, which this research seeks to address.

The evolution of **3D facial animation** has been particularly shaped by the integration of **Generative Adversarial Networks (GANs)**, **Variational Autoencoders (VAEs)**, and more recently, **Transformer-based architectures**. This section outlines the foundational technologies and methods that have contributed to the development of the current state-of-the-art in AI-generated animation.

Copyright to IJARSCT www.ijarsct.co.in



DOI: 10.48175/568





International Journal of Advanced Research in Science, Communication and Technology

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Volume 5, Issue 1, July 2025



2.1 Generative Adversarial Networks (GANs)

Introduced by **Goodfellow et al.**, GANs have revolutionized image synthesis by establishing an adversarial learning framework, where a generator and discriminator are trained simultaneously to produce photorealistic outputs. In facial animation, GANs have been critical for generating high-fidelity facial expressions. Notably, the **First-Order Motion Model** employs keypoint-based motion estimation to animate static images using reference driving videos, enabling flexible and efficient facial motion transfer.

2.2 Variational Autoencoders (VAEs)

VAEs provide a probabilistic framework for learning latent data distributions and generating new instances through sampling. In 3D facial animation, VAEs help capture the variability and subtlety of facial expressions, contributing to smoother frame transitions. Their integration with GANs—as in VAE-GAN architectures—leverages the strengths of both approaches, improving the realism, diversity, and stability of generated animations.

2.3 Transformer-Based Models

Transformers, characterized by their **self-attention mechanisms**, have been adapted for animation tasks to effectively model temporal dependencies in sequential data. Tools like **FaceFormer** utilize transformer networks to synchronize 3D facial animations with speech inputs, offering precise audio-visual alignment. These models excel in capturing long-term dependencies, making them ideal for complex and dynamic animation sequences.

2.4 Emotion-Driven Animation

Incorporating emotional context is vital for producing expressive and believable avatars. **EMOCA** (Emotion Driven Monocular Face Capture and Animation) presents an innovative method by integrating **emotion consistency loss** during training, ensuring the emotional tone of generated animations matches the input imagery. This enhances both the authenticity and expressiveness of facial movements in animated characters.



Copyright to IJARSCT www.ijarsct.co.in



DOI: 10.48175/568





International Journal of Advanced Research in Science, Communication and Technology

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Volume 5, Issue 1, July 2025





III. METHODOLOGY

The proposed system aims to generate realistic 3D facial animations from a single static image using deep learning techniques. The architecture is designed to be user-friendly, efficient, and capable of producing high-quality animations without the need for extensive manual intervention.

3.1. System Overview

The system comprises the following key components:

Input Acquisition: Users provide a static facial image and a driving video that contains the desired facial movements. **Preprocessing**: The inputs undergo preprocessing steps, including face detection, alignment, and normalization, to ensure consistency and compatibility with the model.

Facial Landmark Detection: Utilizing MediaPipe, the system detects and maps 468 facial landmarks, capturing the geometric structure of the face.

Copyright to IJARSCT www.ijarsct.co.in



DOI: 10.48175/568





International Journal of Advanced Research in Science, Communication and Technology

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Volume 5, Issue 1, July 2025



Motion Extraction: The driving video's facial movements are analyzed to extract motion vectors corresponding to the detected landmarks.

Animation Generation: A hybrid VAE-GAN model synthesizes the animated frames by applying the extracted motion vectors to the static image, producing a sequence of frames that depict the desired facial expressions.

Post-Processing: The generated frames are refined using image enhancement techniques to improve visual quality and realism.

Output Compilation: The final frames are compiled into an animation video in formats such as MP4 or GIF.

3.2. Detailed Workflow

Step 1: Input Acquisition

Users upload a high-resolution static image of a face and a short driving video. The system ensures that the inputs meet the required specifications for optimal processing.

Step 2: Preprocessing

Face Detection: OpenCV is employed to detect faces within the inputs. **Alignment**: Detected faces are aligned based on eye positions to standardize orientation. **Normalization**: Pixel values are normalized to facilitate efficient model training and inference.

Step 3: Facial Landmark Detection

MediaPipe's Face Mesh solution detects 468 3D facial landmarks. These landmarks provide a detailed map of facial features, essential for accurate motion transfer

Step 4: Motion Extraction

The system analyzes the driving video to extract motion vectors corresponding to the detected landmarks. Temporal smoothing techniques are applied to ensure consistency across frames.

Step 5: Animation Generation

A VAE-GAN model is trained to generate realistic facial animations.

VAE Component: Encodes the static image into a latent space, capturing essential features.

GAN Component: Decodes the latent representation and applies the motion vectors to generate animated frames.

The model is trained using a combination of reconstruction loss, adversarial loss, and perceptual loss to ensure highquality outputs.

Step 6: Post-Processing

Generated frames undergo enhancement processes, including: Color Correction: Adjusting color balance to match the original image. Sharpness Enhancement: Applying filters to enhance image clarity. Artifact Removal: Eliminating any visual artifacts introduced during generation.

Step 7: Output Compilation

Enhanced frames are compiled into a cohesive animation using FFmpeg. The final output is provided in user-friendly formats such as MP4 or GIF.

Copyright to IJARSCT www.ijarsct.co.in



DOI: 10.48175/568





International Journal of Advanced Research in Science, Communication and Technology

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Volume 5, Issue 1, July 2025



3.3. Visual Illustrations

Figure 1: System Architecture



This diagram illustrates the end-to-end workflow of the proposed system, highlighting each component's role in the animation generation process.





Depicts the 468 facial landmarks detected by MediaPipe, providing a comprehensive map of facial feature

Copyright to IJARSCT www.ijarsct.co.in



DOI: 10.48175/568





International Journal of Advanced Research in Science, Communication and Technology

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Volume 5, Issue 1, July 2025



Figure 3: VAE-GAN Model Architecture



Showcases the integration of VAE and GAN components in the model, facilitating realistic animation generation



Figure 4: Sample Animation Frames

Displays a sequence of frames generated by the system, demonstrating the smooth transition of facial expressions

Copyright to IJARSCT www.ijarsct.co.in



DOI: 10.48175/568





International Journal of Advanced Research in Science, Communication and Technology

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Volume 5, Issue 1, July 2025



IV. RESULTS AND DISCUSSION

The proposed deep learning-based system for 3D facial animation was rigorously evaluated through both **quantitative metrics** and **qualitative assessments**, using a diverse dataset comprising real and synthetic images. The experiments were conducted using **benchmark datasets** such as **VoxCeleb** and **YouTube Faces**, which are widely used in facial recognition and synthesis research.

4.1 Quantitative Evaluation

The system was evaluated using benchmark datasets like **VoxCeleb** and **YouTube Faces**. The input consisted of **real and synthetic facial images** animated using short driver videos. The results were analyzed based on key performance metrics including **SSIM**, **FID**, and **KSR**.

The animation achieved:

SSIM of 0.87, indicating strong structural similarity with the source image.

FID of 23.4, reflecting high-quality and photorealistic generation.

Average FPS of 23, enabling real-time rendering performance.

Keypoint Stability Ratio (KSR) of 91%, ensuring smooth and stable transitions between frames.

These metrics confirm the system's capability to generate structurally accurate, visually realistic, and temporally consistent facial animations at real-time speeds.

Additionally, the system demonstrated efficient performance on standard hardware (e.g., NVIDIA GPUs), with average inference times under **40 ms per frame**, supporting deployment in real-time applications.

4.2 Qualitative Analysis

User experience was captured using **Mean Opinion Scores (MOS)**. Participants rated the generated animations based on realism, identity preservation, and smoothness. The system achieved an average **MOS rating of 4.2 out of 5**, reflecting a high level of satisfaction.

Users particularly noted the system's ability to retain identity features and express emotions convincingly, even under different lighting conditions and facial orientations.

4.3 Versatility and Generalization

The system also supports **diverse face types** — including **avatars, cartoon characters, and statues** — making it highly **versatile**. Its capability to generalize across a variety of facial domains highlights its adaptability to multiple use cases beyond human animation.

Moreover, the system handled varying facial structures, expressions, and skin textures without significant loss of quality, demonstrating the robustness of the **MediaPipe landmark detection** and the hybrid **VAE-GAN architecture**.

4.4 Comparative Performance

Compared to existing frameworks like FOMM (First Order Motion Model) and X2Face, the proposed system delivers:

Better structural preservation, as evidenced by SSIM results.

More photorealistic outputs, indicated by the improved FID.

Higher temporal coherence, ensured by the KSR metric.

Real-time rendering capabilities, which are often limited in other approaches.

4.5 Limitations and Future Scope

Despite its strengths, the system has certain limitations:

Performance slightly degrades with extreme head poses or partial occlusions.

Some minor temporal inconsistencies may appear during rapid or complex facial movements.

Generating higher-resolution animations can increase GPU memory consumption.

Copyright to IJARSCT www.ijarsct.co.in



DOI: 10.48175/568





International Journal of Advanced Research in Science, Communication and Technology

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Volume 5, Issue 1, July 2025



To address these, future work may explore:

Transformer-based architectures for improved temporal modeling.

Self-supervised learning techniques to reduce data dependency.

Integration of audio-driven synchronization to enhance naturalness and interactivity.



Fig. 3 Comparison of SSIM, FPS across Models

V. CONCLUSION

This project introduces an innovative and efficient method for creating 3D facial animations using deep learning. By combining the strengths of GANs and VAEs, the system can bring a still image to life in real time—without the need for expensive motion-capture equipment or manual effort.

With strong results like an SSIM of 0.87, FID of 23.4, and 23 frames per second, the system proves to be both accurate and fast, making it suitable for practical use. Since it's built on a cloud-based platform, even users with basic hardware can access and use it easily.

Looking ahead, the project opens the door for exciting enhancements like **full-body animation**, **web-based interfaces**, and **lip-syncing animations from audio input**. Overall, this system offers a powerful, accessible, and cost-effective solution for generating 3D animations, setting a new standard in the field.

VI. ACKNOWLEDGMENT

We would like to express our heartfelt gratitude to our project guide, **Prof. Pushpa Tandekar**, whose continuous support, insightful feedback, and expert guidance have been instrumental throughout the course of our research. Her patience, motivation, and in-depth knowledge inspired us to push our limits and gave this project a strong foundation.

We are also deeply thankful to the **Department of Computer Science and Engineering** and our college for providing us with the necessary infrastructure, resources, and a positive academic environment that facilitated our development and experimentation.

Our sincere thanks go to all the **faculty members** who offered their valuable suggestions and technical advice whenever needed. Your encouragement and input helped us refine our ideas and troubleshoot challenges at various stages of the project.

We also acknowledge the unwavering support of our **friends and batchmates**, whose constructive feedback, moral support, and collaborative spirit were truly appreciated.

Lastly, we are grateful to our **families**, who stood by us with their constant encouragement, love, and understanding throughout this journey. Their support played a silent yet crucial role in helping us complete this project successfully.

REFERENCES

- [1]. Siarohin et al., 'First Order Motion Model for Image Animation,' CVPR 2019.
- [2]. E. Zakharov et al., 'Few-Shot Adversarial Learning of Talking Heads,' ICCV 2019.

Copyright to IJARSCT www.ijarsct.co.in



DOI: 10.48175/568





International Journal of Advanced Research in Science, Communication and Technology

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Volume 5, Issue 1, July 2025



- [3]. J. Thies et al., 'Deep Video Portraits,' ACM TOG, 2018.
- [4]. Y. Deng et al., 'Accurate 3D Face Reconstruction,' arXiv:1903.08527.
- [5]. M. Abu Zeid et al., 'Real-Time Facial Animation using Deep Learning,' IEEE TVCG 2021.
- [6]. T. Karras et al., 'A Style-Based Generator Architecture for GANs,' CVPR 2019.
- [7]. M. Liang et al., 'Liquid Warping GAN++,' arXiv:2003.04013.
- [8]. K. Olszewski et al., 'High-Fidelity Facial Avatars from a Single Image,' arXiv 2023.
- [9]. Artificial Neural Network, May 2022, DOI: 10.17148/IJARCCE.2022.115196, Conference: International Journal of Advanced Research in Computer and Communication Engineering.
- [10]. python.net, December 2022, DOI:10.17148/IJARCCE.2022.111237, Conference: International Journal of Advanced Research in Computer and Communication Engineering.
- [11]. Combining Vedic & Traditional Mathematic Practices for Enhancing Computational Speed in Day-To-Day Scenarios, Speed in Day-To-Day Scenarios, Conference: Industrial Engineering Journal ISSN: 0970-2555 Website: www.ivyscientific.org, At: Industrial Engineering Journal ISSN: 0970-2555, Website: www.ivyscientific.org. (UGC JOURNAL)
- [12]. Research on Techniques for Resolving Big Data Issues ,May 2022,DOI:
- [13]. 10.17148/IJARCCE.2022.115192 ,Conference: International Journal of Advanced Research in Computer and Communication Engineering
- [14]. Photometric and spectroscopic analysis of the Type II SN 2020jfo with a short plateau, November 2022
- [15]. DOI:10.48550/arXiv.2211.02823 ,License CC BY 4.0.
- [16]. Research on Data Mining, May 2022, DOI: 10.17148/IJARCCE.2022.115176, Conference: International Journal of Advanced Research in Computer and Communication Engineering.
- [17]. Research on Association Rule Mining Algorithms , May 2022 , DOI: 10.17148/IJARCCE.2022.115152 ,Conference: International Journal of Advanced Research in Computer and Communication Engineering.
- [18]. STUDY on INTERNET of THINGS BASED APPLICATION ,May 2022 , DOI: 10.17148/IJARCCE.2022.115179 , Conference: International Journal of Advanced Research in Computer and Communication Engineering.
- [19]. An Efficient Way to Detect the Duplicate Data in Cloud by using TRE Mechanism , May 2022 ,DOI:10.17148/IJARCCE.2022.115139 ,Conference: International Journal of Advanced Research in Computer and Communication Engineering , Volume:11.
- [20]. Using Encryption Algorithms in CC for Data Security and Privacy , May 2022 , DOI:10.17148/IJARCCE.2022.115149



