Text-To-Image Generation Using Deep Learning

Kunduru Chaitanya Reddy¹, Mrs. A. Tulasi²

¹Student, Roll No: 323206454008, Andhra University College of Engineering, Visakhapatnam, Andhra Pradesh

²Assistant Professor, Department of Computer Science and Systems Engineering, Andhra University College of Engineering, Visakhapatnam, Andhra Pradesh

Abstract—This project explores the fascinating intersection of Natural Language Processing (NLP) and deep learning to create a text-to-image generation system. Leveraging the power of Stable Diffusion models, this application translates textual descriptions into corresponding visual representations. The project utilizes the diffusers and transformers libraries within a Python environment, optimized for execution on Google Colab's T4 GPU. A key focus is on user interaction, facilitated through a Gradio-based web interface. This interface allows users to input text prompts, select from various pre-trained Stable Diffusion models, and finetune generation parameters like inference steps and guidance scale. Furthermore, the system incorporates an image-to-image editing feature, enabling users to upload existing images and modify them based on textual prompts, controlled by a strength parameter. The project addresses the computational demands of deep learning by employing half-precision floating-point operations (float16) for memory efficiency on the GPU. This project demonstrates the practical application of deep learning for creative content generation and highlights the importance of user-friendly interfaces in democratizing access to advanced AI technologies.

Index Terms—Artificial Intelligence, Deep Learning, Text-to-Image Generation, Stable Diffusion, Natural Language Processing, Gradio, Image-to-Image Editing.

I. INTRODUCTION

The convergence of Natural Language Processing (NLP) and deep learning has opened up exciting new avenues in creative content generation. Among these, the ability to translate textual descriptions into visual representations stands out as a particularly compelling application. This project focuses on developing a robust and user-friendly text-to-image generation system, leveraging the power of Stable Diffusion models. This technology represents a significant

advancement in generative AI, enabling the creation of high-quality, diverse, and contextually relevant images from textual prompts. This capability has profound implications across various fields, from art and design to education and entertainment.

The core objective is to build a practical application demonstrating the potential of Stable Diffusion for creative image synthesis. The system is designed to be accessible and efficient, utilizing resources like Google Colab and its T4 GPU. We employ Python as the primary programming language, along with key deep learning libraries such as diffusers and transformers. The diffusers library simplifies loading, configuring, and running these complex models, while transformers provides access to pre-trained tokenizers and NLP components essential for processing textual input.

A crucial aspect is the emphasis on user experience via a Gradio-based web interface. This interface allows users to input text prompts, select from a range of available Stable Diffusion models, and adjust key parameters that influence the generation process. These parameters include the number of inference steps and the guidance scale. Beyond basic generation, an image-to-image editing feature is incorporated, allowing users to upload an existing image and guide modifications using textual prompts, controlled by a strength parameter. Given the computational demands, resource optimization is managed by utilizing half-precision floating-point operations (float16) to reduce memory consumption and accelerate computations on the GPU.

II. RELATED WORK/LITERATURE SURVEY

The proposed system is built upon foundational concepts and advancements in deep learning for

© November 2025 | IJIRT | Volume 12 Issue 6 | ISSN: 2349-6002

generative AI:

- Stable Diffusion (Latent Diffusion Model LDM): This model operates in alower-dimensional latent space, allowing for efficient training and sampling, which significantly reduces computational requirements and makes large-scale image generation feasible.
- DALL-E 2 (Transformer-based diffusion model): Employs a transformer architecture, which is effective in capturing complex relationships between text and images, leading to highly creative and contextually relevant generation.
- Imagen (Cascaded Diffusion Models): Focuses on generating high-resolution images by using a cascaded diffusion model approach, progressively refining the image from low to high resolution.
- VQ-GAN+CLIP: Combines Vector Quantized GANs (VQ-GANs) for image generation with CLIP (Contrastive Language-Image Pre-training) to align the generated images with their textual descriptions, ensuring consistency.
- DALL-E: A pioneering transformer-based model that demonstrated the feasibility and potential of using transformers for text-to-image synthesis, paving the way for subsequent research.
- GLIDE (Diffusion Model): Emphasizes the importance of effective text conditioning in guiding the image generation process within a diffusion framework.
- AttnGAN (GAN with Attention Mechanism):
 Explored using attention mechanisms to enable finer-grained control by allowing the model to selectively focus on relevant parts of the text input when generating different parts of the image.

The efficiency of latent diffusion models (like Stable Diffusion) and the text-image relationship capture via transformers (like DALL-E 2) directly inform the design of the proposed system.

III. PROBLEM STATEMENT & OBJECTIVES

A. Problem Statement

The increasing accessibility of advanced text-to-image generation models is hindered by several challenges. Current solutions often possess steep learning curves, requiring specialized technical expertise to operate effectively, thus limiting their use to technically

proficient individuals. Furthermore, the computational demands of these models pose a significant barrier for users with limited access to powerful hardware or constrained cloud computing resources, making effective resource management in environments like Google Colab crucial. Additionally, users face difficulties navigating the available models, understanding the impact of various generation parameters, and incorporating their own images for editing.

B. Objectives

This project addresses these challenges by proposing a user-friendly and resource-conscious text-to-image generation system. The specific objectives are:

- Develop a user-friendly interface for text-toimage generation, enabling intuitive prompt input and parameter adjustment.
- Implement a system for selecting and loading various pre-trained Stable Diffusion models, providing users with diverse artistic styles.
- Incorporate an image-to-image editing feature, allowing users to modify existing images based on textual descriptions.
- Optimize the system for resource-constrained environments like Google Colab, utilizing techniques such as half-precision floating-point operations (float16).
- Ensure robust and reliable image generation, addressing potential challenges related to memory management and library compatibility.

IV. PROPOSED SYSTEM ARCHITECTURE

The proposed system adopts a modular architecture designed for efficiency and user-friendliness, drawing inspiration from latent diffusion models and transformer encoding approaches.

A. High-level components

- User Interface: A Gradio-based web interface serves as the user-friendly interaction point. It allows users to input text prompts, select models, adjust parameters (e.g., guidance scale, inference steps), upload images (for editing), and view generated images.
- Model Management: This module handles the loading and switching of pre-trained Stable Diffusion models. It manages a dictionary of

© November 2025 | IJIRT | Volume 12 Issue 6 | ISSN: 2349-6002

- available models and efficiently loads the selected model onto the designated device (GPU or CPU).
- Image Generation/Editing Engine: This core module uses the diffusers library to execute the generation task. It handles generating images from text prompts and performing image editing based on user input.
- Resource Management: This component optimizes resource usage, particularly memory management in Google Colab. It utilizes halfprecision floating-point operations (float16) and manages the clearing of CUDA cache to prevent memory issues.

B. System Workflow

- The User provides a text prompt (and optionally uploads an image) through the Gradio interface.
- The User selects a model and adjusts parameters (e.g., strength parameter for editing).
- The Model Management module loads the chosen model.
- The transformers library is employed to tokenize and encode the text prompt into numerical representations suitable for the Stable Diffusion model.
- The Image Generation/Editing Engine processes the encoded prompt (and input image, if provided) using the selected model. For editing, a strength parameter controls the degree of alteration to the original image.
- The generated image is displayed in the Gradio interface.

V. METHODOLOGY & IMPLEMENTATION DETAILS

A. Tech stack & rationale

- Models: Pre-trained Stable Diffusion models are utilized, accessible via the diffusers library. These models, trained on extensive image-text pairs, provide diverse and contextually appropriate generation capabilities.
- Libraries: diffusers simplifies loading, configuration, and execution of the diffusion pipelines. transformers handles the necessary text processing, specifically tokenization and encoding of the input prompt.
- User Interface: Gradio is the foundation of the

- front-end, providing a dynamic and simple interface.
- Resource Optimization: Half-precision floatingpoint operations (float16) are implemented in the model loading and execution process to ensure memory efficiency on the Google Colab GPU.

B. Key Functionalities

- Model Selection: The system offers a choice of several Stable Diffusion models via a dynamic dropdown menu, allowing users to experiment with different artistic styles and characteristics from the same text prompt.
- Text Processing (Feature Extraction): The input text prompt is tokenized and encoded using the transformers library, converting the text into numerical representations that capture its semantic meaning for the Stable Diffusion model.
- Image-to-Image Editing: This feature is crucial, allowing users to upload an image and use a text prompt to guide modifications. The model generates a modified image, with the strength parameter controlling the blend between the original image and the text-guided generation.

C. Dataset

The project utilizes the existing training of the pretrained Stable Diffusion models. These models are readily available through the diffusers library. No additional training of the Stable Diffusion models or custom dataset collection/splitting is performed within this project.

VI. RESULTS AND FUTURE WORK

The primary outcome of this project is a functional and user-friendly text-to-image generation system deployed on Google Colab. This system successfully bridges the gap between sophisticated deep learning models and non-technical users by providing an intuitive interface for prompt input, model selection, and parameter adjustment. The system's optimization via float 16 operations ensures effective use of readily available platforms.

The key results demonstrated include:

• Successful implementation of a Stable Diffusion pipeline using diffusers and transformers.

© November 2025 | IJIRT | Volume 12 Issue 6 | ISSN: 2349-6002

- User-centric design through the Gradio interface, making complex AI technology accessible.
- Resource efficiency achieved through halfprecision operations, enabling robust generation in a resource-constrained cloud environment.
- Expanded creative potential via the image-toimage editing feature and diverse model selection.

A. Future Work

Future work may include:

- Integrating more advanced prompt engineering techniques to enhance control over the output.
- Implementing additional image editing functionalities, such as inpainting or outpainting.
- Exploring further optimization techniques for resource efficiency and faster generation speed.
- Conducting user studies to evaluate the system's usability and creative effectiveness in detail.

VII. CONCLUSION

We presented a user-friendly and resource-conscious text-to-image generation system leveraging Stable Diffusion models and implemented with the diffusers and transformers libraries. The project achieved its objective of creating a practical application that efficiently translates textual descriptions into visual realities. Key features, including the Gradio interface, flexible model selection, and the image-to-image editing functionality, successfully democratize access to advanced generative AI. The utilization of half-precision floating-point operations ensures effective resource management, proving the viability of deploying such complex models on readily available cloud platforms.

REFERENCES

- [1] Rombach, R., Blattmann, A., Lorenz, D., Esser, P., & Ommer, B. "High-resolution image synthesis with latent diffusion models." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2022.
- [2] Yu, J., Li, W., Bao, J., Yang, Y., & Song, S. "Vector quantize and classify: A simple approach for textto-image synthesis." Proceedings of the IEEE/CVF International Conference on Computer Vision, 2021.

- [3] Ramesh, A., Pavlov, M., Goh, G., Gray, S., Voss, C., Radford, A., & Sutskever, I. "Zero-shot text-toimage generation." International Conference on Machine Learning, PMLR, 2021.
- [4] Karras, T., Laine, S., & Aila, T. "A style-based generator architecture for generative adversarial networks." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2019.
- [5] Xu, T., Zhang, P., Huang, Q., Zhang, H., Gan, Z., Huang, X., & He, X. "Attngan: Fine-grained text to image generation with attentional generative adversarial networks." Proceedings of the IEEE conference on computer vision and pattern recognition, 2018.
- [6] Zhang, H., Xu, T., Li, H., Zhang, S., Wang, X., Huang, X., & Metaxas, D. N. "Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks." Proceedings of the IEEE international conference on computer vision, 2017.