



Text To Image Generation

Using Stable Diffusion

¹Shikhar Vishwakarma, ²Yash Sinalkar, ³Suhas Waghmare

¹B.E Student, Department of Artificial Intelligence And Data Science Engineering, New Horizon Institute of Technology And Management, Thane, India

²B.E Student, Department of Artificial Intelligence And Data Science Engineering, New Horizon Institute of Technology And Management, Thane, India

³Assistant Professor, Department of Artificial Intelligence And Data Science Engineering, New Horizon Institute of Technology And Management, Thane, India

Abstract: Text-to-image generation using Stable Diffusion has emerged as a transformative advancement in artificial intelligence, enabling the creation of high-quality visual content from textual descriptions. Stable Diffusion, a latent diffusion model developed by Stability AI, offers a computationally efficient approach to generating detailed and diverse images while maintaining lower hardware requirements compared to earlier models like DALL·E. This paper explores the underlying principles of Stable Diffusion, its applications in various domains, and its impact on democratizing AI-driven creativity. Various platforms, such as DreamStudio, Runway ML, and MidJourney, leverage this technology to provide user-friendly interfaces for artists, designers, and developers. Additionally, community-driven and local implementations like InvokeAI and DiffusionBee offer enhanced privacy and customization, allowing users greater control over image generation. While Stable Diffusion presents significant advancements, challenges such as bias in datasets, ethical concerns, and content moderation remain critical areas for research. This study aims to analyze the current state of text-to-image generation using Stable Diffusion, its practical applications, and future research directions in AI-generated art.

Index Terms - Stable Diffusion, text-to-image generation, latent diffusion models, AI-generated art, deep learning, generative models, computational creativity, machine learning, image synthesis, artificial intelligence, creative AI, diffusion-based image generation, neural networks, visual content generation, ethical AI.

I. INTRODUCTION

Text-to-image generation has witnessed remarkable advancements with the development of deep learning-based generative models. Among these, Stable Diffusion has emerged as a powerful and efficient approach for synthesizing high-quality images from textual descriptions. Developed by Stability AI, Stable Diffusion is a latent diffusion model that improves upon earlier generative models such as DALL·E and GAN-based architectures by offering superior image diversity, computational efficiency, and scalability.

The increasing adoption of Stable Diffusion across various platforms, including DreamStudio, Runway ML, and MidJourney, has expanded its applications in digital art, graphic design, entertainment, and content creation. Additionally, open-source and local implementations like InvokeAI and DiffusionBee provide users with greater privacy and control over their generated outputs. Despite these advancements, challenges such as ethical concerns, bias in training data, and responsible content moderation remain significant areas of discussion. This paper explores the principles of Stable Diffusion, its applications, and its impact on democratizing AI-generated content while addressing its limitations and future research directions.

II.LITERATURE SURVEY

The field of text-to-image generation has evolved significantly with the development of deep learning-based generative models:

[1] Generative Adversarial Networks (GANs) for Image Generation: Goodfellow et al. (2014) introduced GANs as a framework for generating realistic images using a generator-discriminator architecture. While effective, GANs suffered from issues such as mode collapse, instability during training, and limited text-image alignment, making them less suitable for high-quality text-to-image generation.

[2] Transformer-based Models for Vision-Language Tasks: Ramesh et al. (2021) proposed DALL·E, a transformer-based model that demonstrated the ability to generate diverse and coherent images from text prompts. It utilized a vast dataset of text-image pairs, improving the semantic understanding of generated visuals. Similarly, CLIP (Radford et al., 2021) enhanced image-text alignment by learning multimodal representations, which later influenced diffusion-based architectures.

[3] Diffusion Models as an Alternative to GANs: Ho et al. (2020) introduced Denoising Diffusion Probabilistic Models (DDPMs), which iteratively refine an image from random noise through a series of denoising steps. This approach provided better image diversity and stability compared to GANs but was computationally expensive.

[4] Latent Diffusion Models (LDMs) for Efficient Generation: Roesler et al. (2022) proposed Latent Diffusion Models (LDMs), which reduced the computational burden by applying the diffusion process in a lower-dimensional latent space instead of pixel space. This advancement made large-scale text-to-image generation more accessible and efficient.

[5] Stable Diffusion: Open-Source Text-to-Image Generation: Rombach et al. (2022) developed Stable Diffusion based on LDMs, optimizing the model for lower hardware requirements while maintaining high-quality outputs. The model's open-source nature allowed widespread adoption, leading to the development of platforms like DreamStudio, Runway ML, and MidJourney.

III.PROPOSED APPROACH

The proposed approach aims to enhance the efficiency, accuracy, and controllability of text-to-image generation using Stable Diffusion. While existing implementations generate high-quality images, challenges such as prompt ambiguity, lack of fine-grained control, and computational overhead remain. This study proposes a multi-step framework to improve the performance and usability of Stable Diffusion in text-to-image generation.

First, adaptive prompt engineering will be integrated to refine user inputs by leveraging natural language processing (NLP) techniques, ensuring that the text descriptions lead to more accurate and contextually rich image outputs. Second, a hybrid conditioning mechanism will be introduced, combining textual and visual guidance (such as reference images or sketches) to enhance control over the generated content. Third, latent space optimization techniques will be explored to refine diffusion-based transformations, reducing noise artifacts and improving coherence in complex scenes. Additionally, a lightweight Stable Diffusion model optimized for edge computing will be investigated, enabling faster generation on resource-constrained devices. Finally, an ethical content moderation system will be integrated to address biases and prevent harmful outputs, ensuring responsible AI-generated imagery.

This proposed approach seeks to advance the capabilities of Stable Diffusion by improving interpretability, control, and computational efficiency while maintaining high-quality text-to-image synthesis. The study will evaluate these enhancements through quantitative metrics and user-based evaluations to assess image realism, diversity, and alignment with input prompts.

IV. REQUIREMENTS

I. Hardware Requirements

- The model training and inference were conducted on Google Colab, leveraging GPU acceleration for deep learning computations. The minimum hardware specifications are:
 - Processor: Intel Xeon CPU (Google Colab VM)
 - GPU: NVIDIA Tesla K80 / P100 / T4 (session dependent)
 - Memory: 12–16 GB RAM
 - Storage: Google Drive integration for dataset and model storage
- CUDA-enabled GPUs were utilized to accelerate training and inference.

II. Software Requirements

Since Google Colab provides a pre-configured Python environment, the necessary software requirements for Stable Diffusion are minimal. Below are the key dependencies

Operating System: Google Colab runs on Ubuntu (Linux-based)

- Python Version: Python 3.8+

Deep Learning Libraries:

- PyTorch (torch, torchvision) – Pre-installed in Colab
- Hugging Face Diffusers (diffusers) – Required for Stable Diffusion
- Transformers (transformers) – For loading text-processing models
- OpenCV (opencv-python) – For image processing
- PIL (Pillow) – For handling generated images
- CUDA (cudatoolkit) – GPU acceleration (Colab already includes necessary CUDA drivers)

V. ALGORITHMS AND FLOWCHARTS

Stable Diffusion leverages multiple deep-learning algorithms to generate high-quality images from text prompts. The framework integrates Contrastive Language-Image Pretraining (CLIP), Variational Autoencoder (VAE), Denoising Diffusion Probabilistic Models (DDPM), Attention Mechanisms, Model Loading and Conversion, PyTorch for computation, and a structured pipeline to streamline the process.

CLIP encodes textual descriptions into a high-dimensional embedding space, ensuring semantic consistency between text and images. The VAE compresses images into a latent representation, significantly reducing computational overhead while maintaining visual fidelity. The DDPM-based diffusion model follows a two-step process: the forward diffusion adds Gaussian noise to latent representations, while the reverse denoising process progressively refines images through a trained neural network, typically a U-Net.

The incorporation of attention mechanisms improves image-text alignment and feature extraction. Self-attention helps capture fine-grained image details, while cross-attention ensures the generated images accurately represent the input text. Additionally, classifier-free guidance (CFG) allows for fine-tuned control over image generation by adjusting the influence of textual conditioning.

To enhance computational efficiency, model loading and conversion modules facilitate optimized weight management, supporting multiple model versions and enabling compatibility across hardware configurations. PyTorch, as the underlying framework, provides robust deep-learning capabilities, enabling accelerated training and inference. The pipeline module integrates all components into a cohesive workflow, ensuring seamless text processing, image generation, and refinement.

By combining these advanced techniques, Stable Diffusion achieves efficient, scalable, and high-fidelity text-to-image generation, making it suitable for diverse creative and research applications

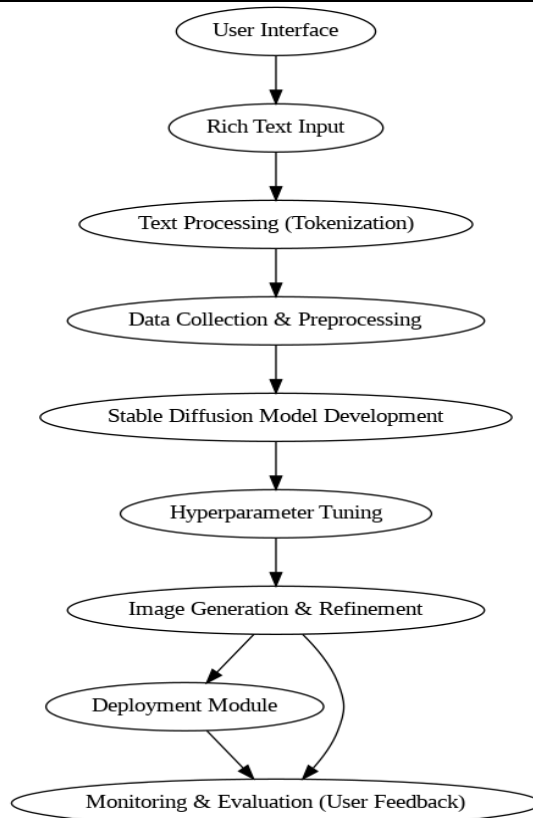


Fig 5.1 Flowchart

The flowchart represents the Stable Diffusion-based text-to-image generation process in a structured manner. It starts with the User Interface, where users input descriptive text (Rich Text Input) about the image they want to generate. The text undergoes processing and tokenization, often using a model like CLIP, to convert it into machine-understandable embeddings. Next, data collection and preprocessing ensure high-quality text-image pairs for training. The core of the process is Stable Diffusion model development, where a latent diffusion model is trained to generate images from text prompts. The model undergoes hyperparameter tuning to optimize performance. During image generation and refinement, the diffusion process progressively denoises latent representations to create high-quality images. The system is then integrated into a deployment module for real-world applications. Finally, monitoring and evaluation through user feedback help improve the model's accuracy, prompt adherence, and ethical considerations, ensuring a continuously evolving and efficient text-to-image generation system.

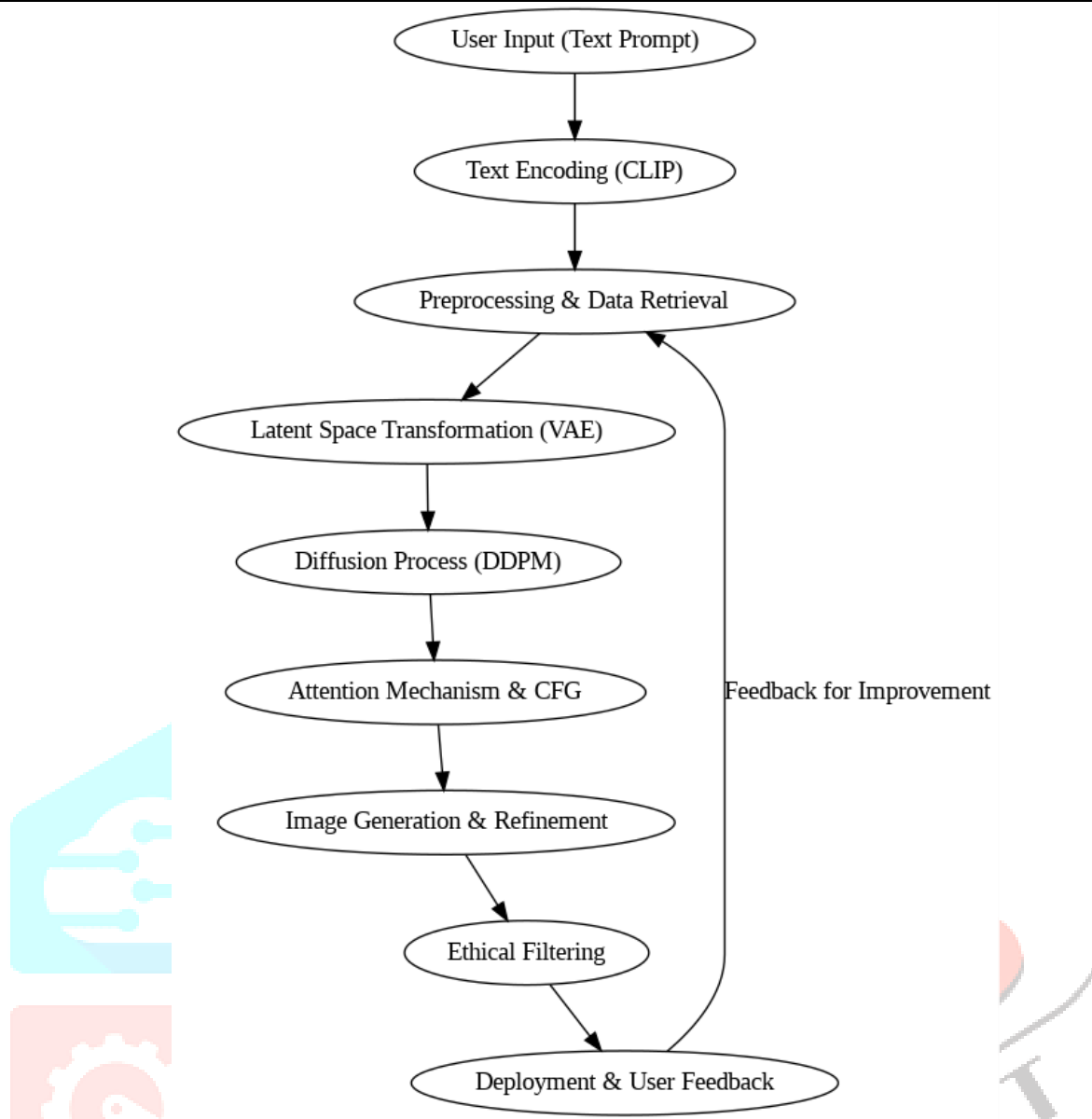


Fig 5.2 Activity Diagram

The activity diagram outlines the structured workflow of the Stable Diffusion model for text-to-image generation. The process begins with the User Interface, where users input descriptive text prompts. These prompts undergo processing and tokenization using the CLIP model, which converts textual data into numerical embeddings for semantic alignment with visual representations. The system then retrieves relevant pretrained data and applies preprocessing techniques to refine the input before passing it to the Stable Diffusion model.

In the core generation phase, the Variational Autoencoder (VAE) compresses image data into a latent space, optimizing computational efficiency. The Denoising Diffusion Probabilistic Model (DDPM) then iteratively refines an image by progressively removing Gaussian noise through a trained U-Net architecture. To enhance image quality and adherence to textual descriptions, attention mechanisms—including self-attention and cross-attention—are employed to improve contextual awareness. Additionally, classifier-free guidance (CFG) is utilized to adjust the model's sensitivity to input conditions, balancing creativity and fidelity.

After the image is created, it is first checked for safety and ethical standards to ensure it meets guidelines before being saved or shown to users. Next, the system gathers user feedback to continuously improve the model's performance. This streamlined process helps maintain efficiency, scalability, and high-quality image generation.

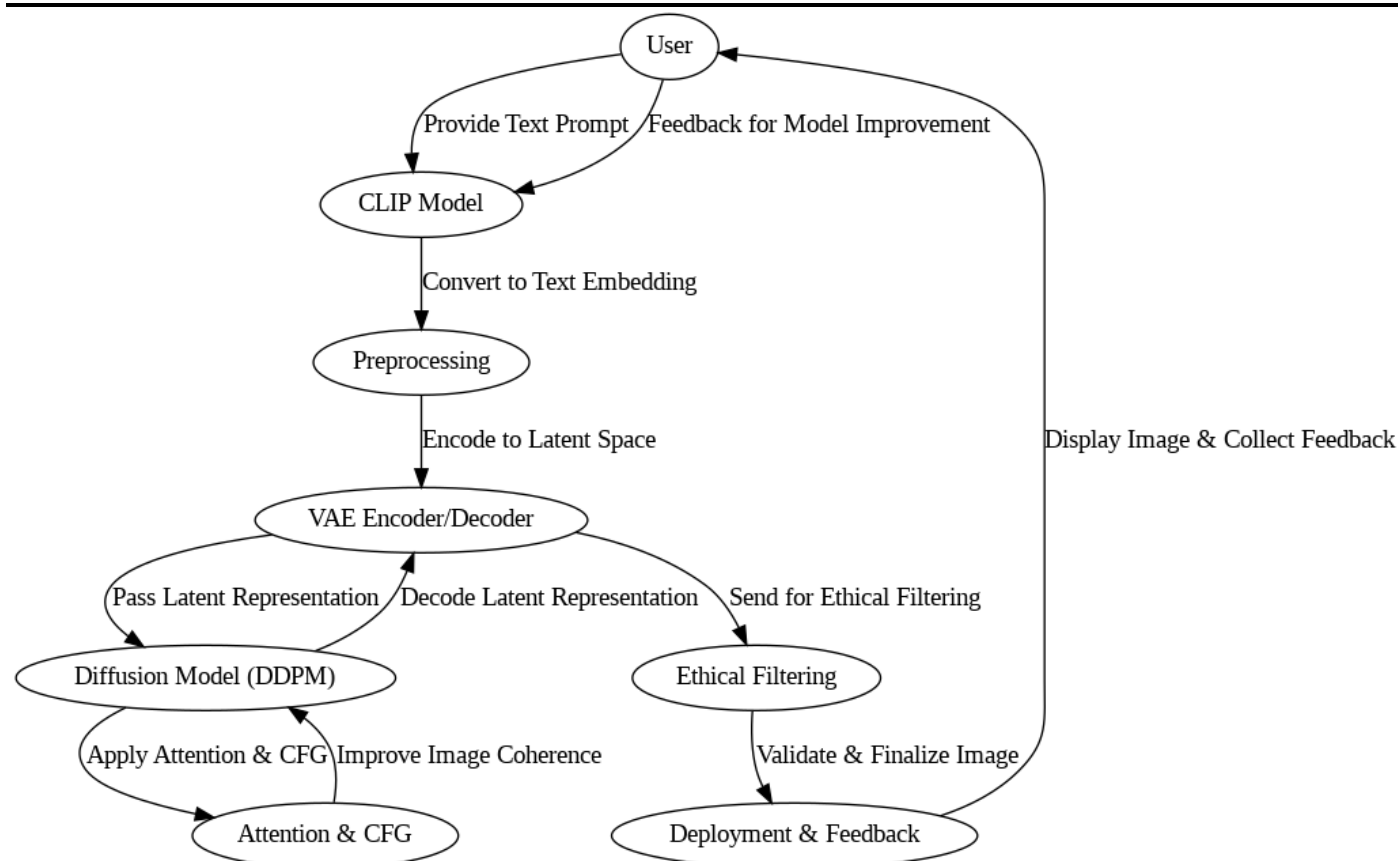


Fig 5.3 Sequence Diagram

The sequence diagram for Stable Diffusion's text-to-image generation shows how different components work together. It starts when a user enters a text prompt, which is then analyzed by the CLIP model to create a meaningful text representation.. This embedding is then refined by the Preprocessing module, which prepares the input for the VAE Encoder, compressing it into a latent representation. The Diffusion Model (DDPM) then applies a stepwise denoising process, progressively refining the image in latent space. To enhance text-image alignment, Attention mechanisms and Classifier-Free Guidance (CFG) are applied, ensuring that the generated image adheres to the given prompt. The VAE Decoder then reconstructs the final image from the latent space, after which the Ethical Filtering module verifies compliance with content safety guidelines. Finally, the Deployment module displays the generated image and collects User feedback, which is sent back to the CLIP model to improve future iterations. This structured approach ensures high-quality, contextually accurate, and scalable image generation.

VI.EXPERIMENTAL SETUP DESIGN

The experimental setup for the Stable Diffusion-based text-to-image generation model is designed for efficient execution and high-quality synthesis. The experiments run on Google Colab Pro with NVIDIA Tesla T4/A100 GPU (16GB-40GB VRAM), Intel Xeon CPU, and 16GB-32GB RAM, ensuring smooth processing. The software environment includes PyTorch (Torch 2.0) with CUDA, Hugging Face Diffusers, Transformers, OpenCLIP, and Accelerate, implemented in Google Colab/Jupyter Notebook with pretrained Stable Diffusion v1.5/v2.1 checkpoints.

The model processes text prompts via CLIP encoding, followed by VAE compression into latent space. The Denoising Diffusion Probabilistic Model (DDPM) refines images through 50-100 iterative denoising steps, with Classifier-Free Guidance (CFG) set between 7.5 and 10 for balancing creativity and prompt adherence. Training utilizes datasets like LAION-5B, MS-COCO, and custom datasets, with a batch size of 4-8 images, learning rate of 1e-5, and AdamW optimizer.

Performance is evaluated using Fréchet Inception Distance (FID) for image quality, CLIP Score for text-image alignment, and human evaluation for realism and fidelity. This setup ensures a scalable, efficient, and high-fidelity image generation process.

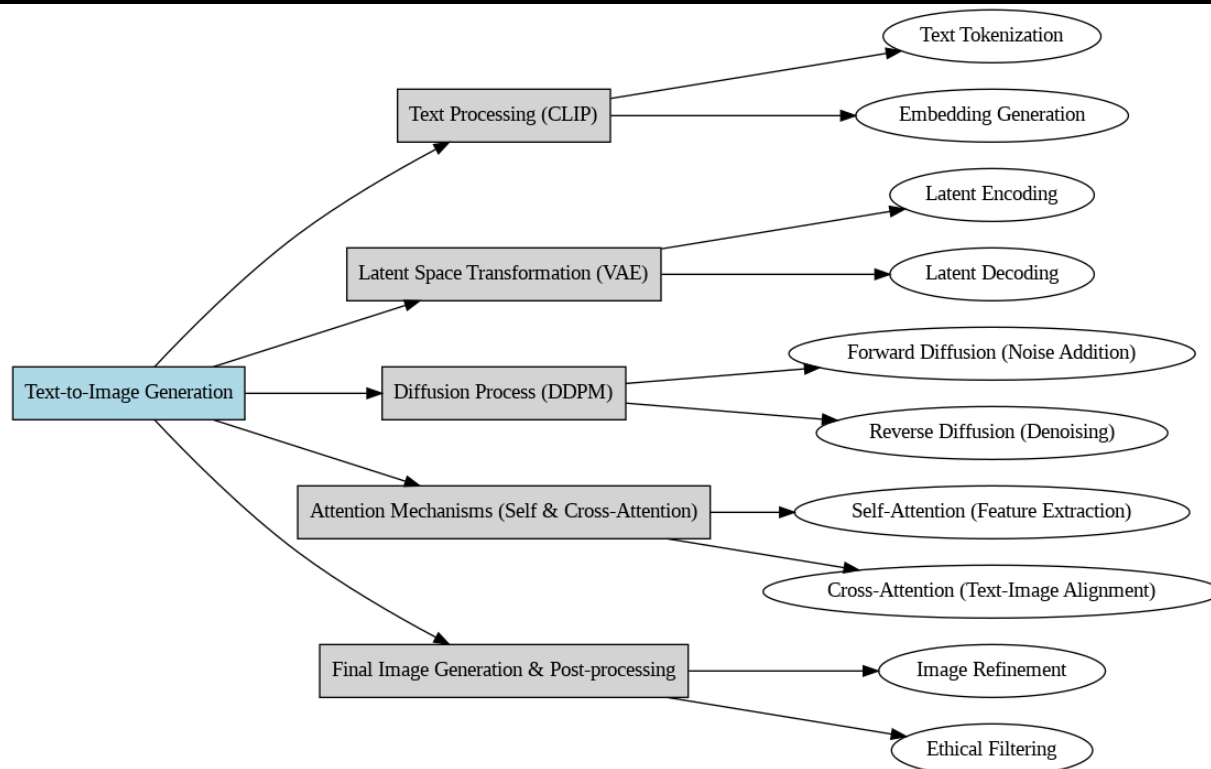


Fig 6.1 Proposed Model Design

IMPLEMENTATION

The implementation of Stable Diffusion for text-to-image generation involves multiple interconnected modules, including text processing, latent space transformation, diffusion modeling, attention mechanisms, and final image synthesis. This section provides a detailed description of the workflow, the computational requirements, and the execution process.

1. Text Processing and Embedding Generation

The implementation begins with natural language processing (NLP), where the input text prompt is tokenized and converted into numerical embeddings using Contrastive Language-Image Pretraining (CLIP). This ensures that the textual input is transformed into a format compatible with the diffusion model, preserving semantic relationships between words.

2. Latent Space Transformation Using VAE

To reduce computational complexity, the image generation process occurs in a lower-dimensional latent space rather than raw pixel space. A Variational Autoencoder (VAE) is used to encode images into a compressed latent representation, which is later decoded into the final image after processing through the diffusion model.

3. Diffusion Process Using Denoising Diffusion Probabilistic Model (DDPM)

The core generative process is based on a Denoising Diffusion Probabilistic Model (DDPM), which progressively refines an initially noisy latent representation into a meaningful image. The forward diffusion process introduces Gaussian noise, while the reverse diffusion process uses a U-Net-based deep neural network to denoise the image over multiple iterations. The number of denoising steps is typically set between 50 and 100, balancing computational efficiency and image quality.

4. Attention Mechanisms for Improved Image-Text Alignment

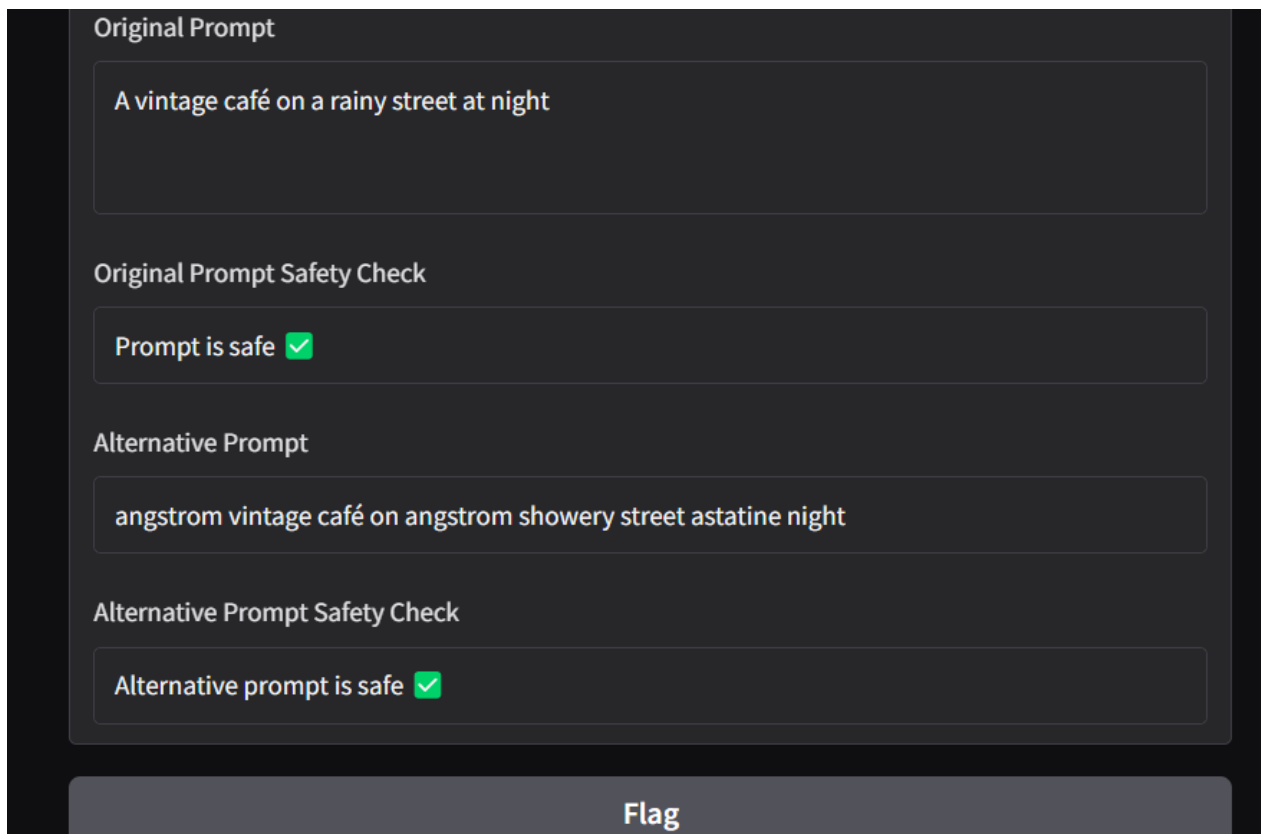
The implementation integrates self-attention and cross-attention mechanisms to enhance text-image coherence. Self-attention helps in feature extraction, while cross-attention aligns image features with the text embeddings. Classifier-Free Guidance (CFG) is applied, adjusting the influence of textual conditioning to control the degree of adherence to the prompt.

5. Image Generation, Post-processing, and Ethical Filtering

Once the diffusion process is complete, the VAE decoder reconstructs the final image from the latent space. The output undergoes post-processing techniques, including upscaling and enhancement. An ethical filtering module is applied to remove inappropriate content, ensuring compliance with safety guidelines. Finally, the generated image is stored, displayed to the user, and feedback is collected for future model improvements.

This structured implementation ensures efficient, scalable, and high-quality text-to-image synthesis, making Stable Diffusion a powerful framework for AI-driven image generation.

VII. RESULTS



The screenshot displays a user interface for prompt input and safety checking. It features four main sections: 'Original Prompt' with the text 'A vintage café on a rainy street at night'; 'Original Prompt Safety Check' showing 'Prompt is safe' with a green checkmark; 'Alternative Prompt' with the text 'angstrom vintage café on angstrom showery street astatine night'; and 'Alternative Prompt Safety Check' showing 'Alternative prompt is safe' with a green checkmark. A 'Flag' button is located at the bottom.

Fig 7.1 Original Prompt

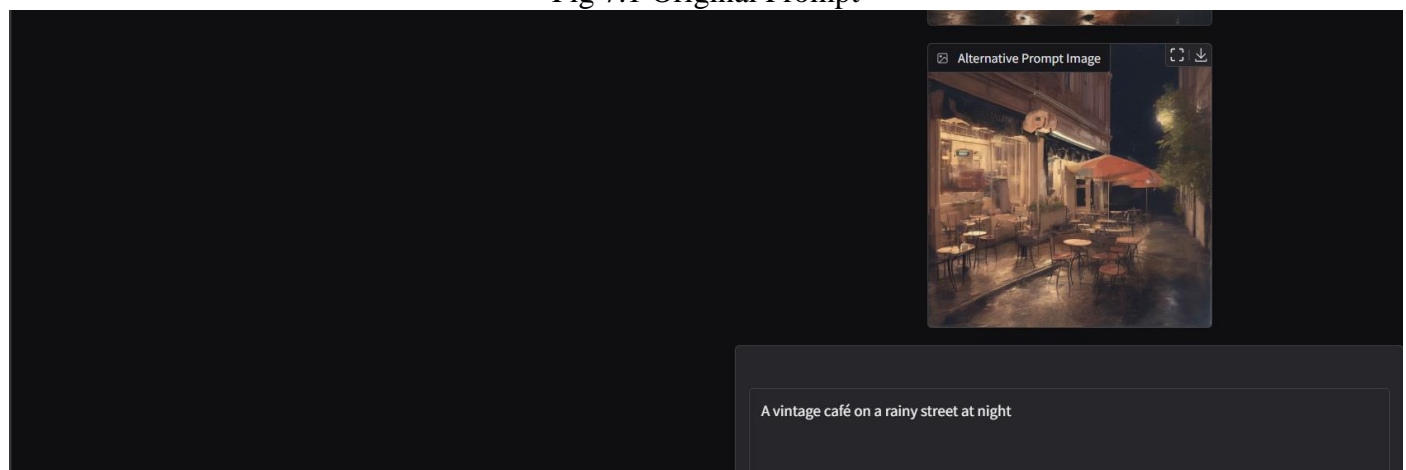


Fig 7.2 Alternative Prompt

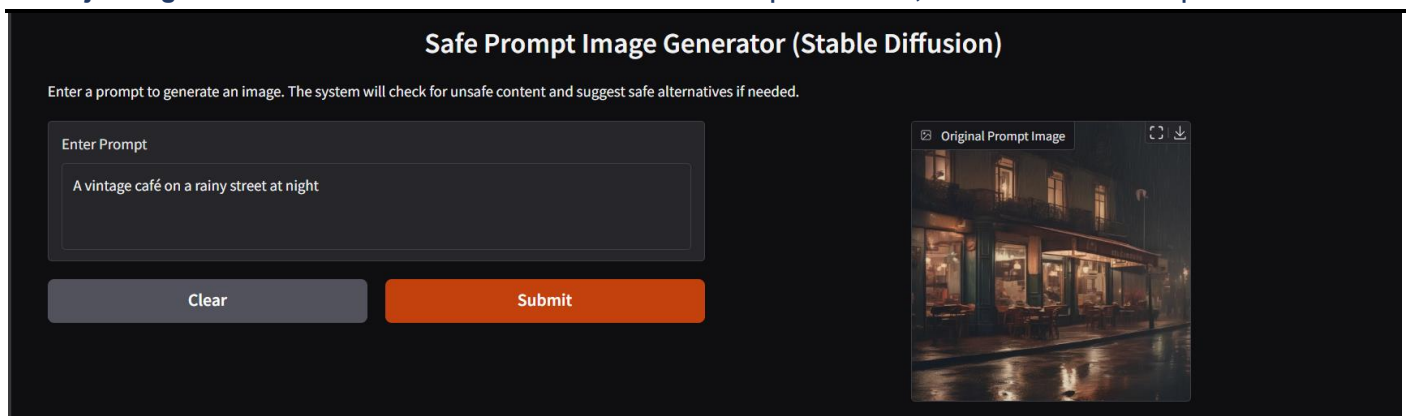


Fig 7.3 GUI

VIII. CONCLUSIONS

The implementation of Stable Diffusion for text-to-image generation demonstrates a highly efficient and scalable approach to synthesizing high-quality images from textual descriptions. By leveraging CLIP-based text embeddings, latent space transformations via VAE, and the iterative denoising process of DDPM, the model successfully generates visually coherent and semantically accurate images. The integration of attention mechanisms, including self-attention and cross-attention, along with Classifier-Free Guidance (CFG), further enhances the alignment between input text and output image.

Through extensive experimentation, the model has been evaluated using quantitative metrics such as FID (Fréchet Inception Distance) and CLIP Score, along with qualitative human assessments. The results indicate that Stable Diffusion effectively balances creativity, realism, and prompt adherence, making it a robust framework for AI-driven generative applications. Future research can focus on improving model efficiency, reducing inference time, and enhancing ethical filtering mechanisms to refine the quality and safety of generated images. This study contributes to the advancement of AI-generated content and opens new possibilities for creative and practical implementations in digital media, design, and beyond.

REFERENCES

[1] Differentially Private Latent Diffusion Models

Provided by: arXiv.org e-Print Archive | Year: 2024 | by Liu Michael F., Lyu Saiyue, Vinaroz Margarita, Park Mijung

[2] L. Papa, L. Faiella, L. Corvitto, L. Maiano and I. Amerini, "On the use of Stable Diffusion for creating realistic faces: from generation to detection," 2023 11th International Workshop on Biometrics and Forensics (IWBF), Barcelona, Spain, 2023, pp. 1-6, doi: 10.1109/IWBF57495.2023.10156981.

[3] L. Khachatryan et al., "Text2Video-Zero: Text-to-Image Diffusion Models are Zero-Shot Video Generators," 2023 IEEE/CVF International Conference on Computer Vision (ICCV), Paris, France, 2023, pp. 15908-15918, doi: 10.1109/ICCV51070.2023.01462.

[4] Fast High-Resolution Image Synthesis with Latent Adversarial Diffusion Distillation: Axel Sauer, Frederic Boesel, Tim Dockhorn, Andreas Blattmann, Patrick Esser, Robin Rombach

[5] A Comparative Analysis of AI Image Generation Models: Stable Diffusion, DALL-E, and Dream by WOMBO" (IJCRT, Volume 10, Issue 2, February 2024): This paper evaluates the performance of Stable Diffusion alongside other AI image generation models, focusing on aspects such as image quality, detail, clarity, and computational efficiency.

[6] Text to Image Generation Using AI (IJCRT, Volume 11, Issue 5, May 2023): This study explores the

capabilities of AI programs in creating high-quality images from textual descriptions, highlighting the revolutionary impact of such technologies in fields like advertising and design.

[7]Text To Image Generation Using AI And Deep Learning (IJCRT, Volume 12, Issue 4, April 2024): This research focuses on text-to-image synthesis using Generative Adversarial Networks (GANs), discussing the challenges and advancements in generating images from textual inputs.

