

DAYANANDA SAGAR UNIVERSITY

KUDLU GATE, BANGALORE – 560068



**SCHOOL OF
ENGINEERING**

Bachelor of Technology

in

COMPUTER SCIENCE AND ENGINEERING

(ARTIFICIAL INTELLIGENCE AND MACHINE LEARNING)

A Project Report On

Harnessing Stable Diffusion Model for High-Resolution

Text-to-Image Synthesis

By

ASHISH PATIL - ENG21AM0013

ATHARVA T - ENG21AM0014

DIKSHA SINHA - ENG21AM0033

TENZIN LUDUP - ENG22AM3011



Under the supervision of

Dr. Vegi Fernando A

Assistant Professor

Computer Science & Engineering (AI & ML)

DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING

(ARTIFICIAL INTELLIGENCE AND MACHINE LEARNING)

SCHOOL OF ENGINEERING

DAYANANDA SAGAR UNIVERSITY

(2024 – 2025)

DAYANANDA SAGAR UNIVERSITY



**SCHOOL OF
ENGINEERING**



Department of Computer Science & Engineering (ARTIFICIAL INTELLIGENCE AND MACHINE LEARNING)

KUDLU GATE, BANGALORE – 560068

Karnataka, India

CERTIFICATE

This is to certify that the project entitled “**Harnessing Stable Diffusion Model for High-Resolution Text-to-Image Synthesis**” is carried out by **ASHISH PATIL (ENG21AM0013), ATHARVA T (ENG21AM0014), DIKSHA SINHA (ENG21AM0033), TENZIN LUDUP (ENG22AM3011)**, bonafide students of Bachelor of Technology in Computer Science and Engineering at the School of Engineering, Dayananda Sagar University, Bangalore, in complete fulfillment for the award of a degree in Bachelor of Technology in Computer Science and Engineering, during the year **2024 - 2025**.

Dr. Vegi Fernando A

Assistant Professor

Dept. of CSE (AIML)

School of Engineering

Dayananda Sagar University

Dr. Vinutha N

Project Co-ordinator

Dept. of CSE (AIML)

School of Engineering

Dayananda Sagar University

Dr. Jayavrinda Vrindavanam

Professor & Chairperson

Dept. of CSE (AIML)

School of Engineering

Dayananda Sagar University

Signature

Signature

Signature

Name of the Examiners:

Signature with date:

1.....

.....

2.....

.....

3.....

.....

DECLARATION

We, **ASHISH PATIL (ENG21AM0013), ATHARVA T (ENG21AM0014), DIKSHA SINHA (ENG21AM0033), TENZIN LUDUP (ENG22AM3011)**, are students of the eighth semester B.Tech in Computer Science and Engineering (AI & ML) at the School of Engineering, Dayananda Sagar University. We hereby declare that the Major Project titled “**Harnessing Stable Diffusion Model for High-Resolution Text-to-Image Synthesis**” has been carried out by us and submitted in complete fulfillment for the award of a degree in **Bachelor of Technology in Computer Science and Engineering** during the academic year **2024–2025**.

Student:

Signature

Name 1: ASHISH PATIL

USN: ENG21AM0013

Name 2: ATHARVA T

USN: ENG21AM0014

Name 3: DIKSHA SINHA

USN: ENG21AM0033

Name 4: TENZIN LUDUP

USN: ENG22AM3011

Place: Bangalore

Date:

ACKNOWLEDGEMENT

It is a great pleasure for us to acknowledge the assistance and support of many individuals who have been responsible for the successful completion of this project work. First, we take this opportunity to express our sincere gratitude to School of Engineering & Technology, Dayananda Sagar University for providing us with a great opportunity to pursue our Bachelor's degree in this institution.

We would like to thank **Dr. Udaya Kumar Reddy K R, Dean, School of Engineering & Technology, Dayananda Sagar University** for his constant encouragement and expert advice. It is a matter of immense pleasure to express our sincere thanks to **Dr. Jayavrinda Vrindavanam, Department Chairman, Computer Science and Engineering (Artificial Intelligence and Machine Learning), Dayananda Sagar University**, for providing right academic guidance that made our task possible.

We would like to thank our guide **Dr. Vegi Fernando A, Assistant Professor, Dept. of Computer Science and Engineering of Artificial Intelligence and Machine Learning Dayananda Sagar University**, for sparing his valuable time to extend help in every step of our project work, which paved the way for smooth progress and fruitful culmination of the project.

We would like to thank our **Project Coordinator Dr. Vinutha N** as well as all the staff members of Computer Science and Engineering (AIML) for their support. We are also grateful to our family and friends who provided us with every requirement throughout the course. We would like to thank one and all who directly or indirectly helped us in the Project work

Contents

1	INTRODUCTION	1
1.1	Using Stable Diffusion:Using Generative AI to Change India’s Digital Environment	2
1.2	Comparisons with Previous Studies:	3
2	PROBLEM DEFINITION AND OBJECTIVE	4
2.1	Problem Definition:	4
2.2	Novelty of Proposed Approach	6
3	LITERATURE SURVEY	7
4	METHODOLOGY	11
4.1	Understanding and Leveraging Stable Diffusion	11
4.2	Lightweight Optimization Techniques:	13
4.3	Model Training and Fine-Tuning:	14
4.4	Deployment on Edge Devices:	14
4.5	Testing and Evaluation:	15
4.6	Real-World Use Case Demonstrations:	15
4.7	Ethical and Safe Deployment:	15
4.8	Evaluation:	15
5	REQUIREMENTS	17
5.1	Functional Requirements:	17
5.2	Non-Functional Requirements:	18
5.3	Enhanced Modules and Techniques:	18
6	RESULT AND ANALYSIS	23
7	CONCLUSION AND FUTURE SCOPE	29

LIST OF FIGURES

Fig .Number	Figure Description	Page Number
Fig 4.1	UNet Flow Diagram	12
Fig 4.2	Stable Diffusion Architecture	13
Fig 4.3	BLEU Score	16
Fig 5.1	Multi-Modal Cross Attention Transformer Model	19
Fig 5.2	LoRA (Low-Rank Adaptation) Working	20
Fig 5.3	HyperNetwork-Based Prompt Embedding	21
Fig 5.4	Time vs FID Score for Denoisers and combinations.	22
Fig 6.1	Astronaut on mountain	24
Fig 6.2	Loss Function	25
Fig 6.3	Comparison of Stable Diffusion, ChatGPT, Gemini	27
Fig 6.4	FID Scores	28

ABSTRACT

This project's main goal is to decode verbal cues into equivalent high-quality visuals by putting an advanced AI diffusion model into practice. The diffusion model is a class of generative AI techniques based on the iterative denoising process. It begins with a disorganized image and carefully refines it over a number of cycles to get a coherent and realistic final. In addition to being visually appealing, this systematic process allows the production of visuals that closely correspond to the spoken description supplied by the user. The project makes use of pre-trained models, such as Stable Diffusion, which are well-known for their ability to generate images from text. Fine-tuning is done on particular datasets that are suited to various use cases in order to improve the model's accuracy and adaptability. The system can process a wider variety of textual inputs thanks to this customization, ranging from straightforward descriptive statements to intricate and subtle prompts. The application of cutting-edge methods like latent space processing, which effectively reduces computing overhead and compresses and manipulates data, guarantees that the system retains high-quality output. A U-Net design, a kind of neural network that excels at producing images with minute details and global structures, is used in the study. The model effectively focuses on relevant sections of the input text when paired with attention mechanisms, ensuring accurate and intricate visual outputs.

A wide range of creative industries, including media, art, and design, benefit from this technology's enhanced content generation and visualization capabilities. E-commerce can be utilized to produce product images, healthcare to enhance medical imaging, and education to help Students visualize complex concepts. Through an examination of diffusion model architectures and principles, the study establishes a connection between literary imagination and digital visualization. It establishes the foundation for further advancements by demonstrating AI's potential in content creation and human-computer interaction.

Chapter 1

INTRODUCTION

This research investigates how to bridge the gap between human creativity and machine-generated visuals by using AI diffusion models to produce photorealistic images from verbal descriptions. The model reduces computing demands while improving output quality by including sophisticated innovations like LoRA, HyperNetwork embeddings, Mixture of Denoisers, and a Multi-Modal Cross-Attention Transformer. E-commerce, advertising, entertainment, education, and digital design can all benefit from this more approachable and effective version of stable diffusion that can be used on edge devices.

- **Context:** The project is a result of generative AI's quick development, especially diffusion models, which have demonstrated remarkable ability to create high-quality, coherent images from noisy data.
- **Objective:** To create a diffusion model that is both lightweight and effective, utilizing cutting-edge AI algorithms to convert textual inputs into photorealistic visuals while preserving device performance on low-power platforms.
- **Scope:** This includes producing a variety of visual outputs, testing deployment on edge devices such as the Jetson Nano and Raspberry Pi, and implementing Stable Diffusion with upgrades (LoRA, HyperNetworks, Mixture of Denoisers, and Multi-Modal Transformers).
- **Significance:** It offers useful applications in digital art, education, advertising, AR/VR, and mobile apps. It also supports AI research on efficiency and multimodal learning.

- Audience: Generative AI is being used by AI researchers, developers, digital artists, educators, and industries to create and visualize content.
- Vision: To drive future innovation in human-computer creative cooperation by democratizing powerful generative AI by making high-quality picture generation useable, interpretable, and accessible across platforms.

1.1 Using Stable Diffusion: Using Generative AI to Change India's Digital Environment

Stable Diffusion's ability to create photorealistic visuals from text is revolutionizing the digital landscape in India. It lowers expenses while increasing efficiency in media, e-commerce, and education. Lightweight versions enable makers even in isolated locations, increasing the accessibility and inclusivity of advanced AI nationwide.

- Stable Diffusion: A generative AI model called Stable Diffusion uses text instructions to produce lifelike pictures. It enables users to turn ideas into graphics by converting noisy data into clear visualizations using diffusion techniques and deep learning.
- Gen AI and its growing relevance in India: India is embracing AI across industries at a rapid pace. Generative AI helps automate content generation in response to growing digital needs, reducing time and offering intelligent, scalable solutions for a range of sectors.
- Application of Stable Diffusion in Indian Sectors: It aids in the visualization of difficult subjects in education. It is used for product images in e-commerce. AI-generated images are used in healthcare for awareness and training, while entertainment benefits from quicker content generation.
- Enhance Accessibility through Lightweight AI: Stable Diffusion can now operate on low-power systems like Raspberry Pi thanks to LoRA and HyperNetworks, opening up access to cutting-edge AI technologies in remote and resource-constrained places.
- Benefits to India's Digital Ecosystem: It promotes quicker innovation and development in digital India by democratizing creative tools, empowering artists in rural areas, and decreasing manual design labor.
- Challenges and Ethical Considerations: Misuse, skewed results, and false information are among the issues. Widespread adoption is further hampered in rural regions by poor digital literacy and inadequate infrastructure.

1.2 Comparisons with Previous Studies:

- **Performance vs GANs and VAEs:** Compared to GANs and VAEs, Stable Diffusion produces images with greater detail and consistency. Additionally, it provides better training stability by lowering common problems like artifacts.
- **Text-to-Image Accuracy:** With more accurate interpretation of text prompts, this model generates visuals that closely resemble the input's intended meaning and specifics.
- **Resource Efficiency:** Real-time utilization is made possible by Stable Diffusion's ability to operate effectively on edge devices and reduce computational costs through methods like LoRA and optimized inference.
- **Scalability and Accessibility:** Stable Diffusion is appropriate for a variety of applications, including low-resource environments, because it is simpler to implement and adjust than previous models.
- **Real World Applications:** Stable Diffusion is more widely used across industries than previous models, fostering creative applications in fields like healthcare, education, and design, among others.

Conclusion

Stable Diffusion is a significant advancement in generative AI that more effectively produces accurate text interpretation and high-quality visuals. Because of its lightweight design, it may be used with low-power devices, increasing accessibility in far-flung locations. It is broadly applicable in industries such as e-commerce, healthcare, and education and promotes innovation, digital inclusion, and expansion in India's digital ecosystem.

Chapter 2

PROBLEM DEFINITION AND OBJECTIVE

2.1 Problem Definition:

The process of turning text into high-quality, photorealistic graphics requires a lot of work and is frequently unavailable on low-power devices. Many current models have issues with speed, accuracy, or resource requirements, which restricts their use, particularly in settings with limited resources like rural India.

- **Approach:** Create a strong AI diffusion model that can produce photorealistic, high-fidelity visuals straight from text inputs while preserving coherence and catching fine details.
- **Enhance Model Efficiency:** Reduce computing needs and enable real-time inference by incorporating lightweight techniques like Low-Rank Adaptation (LoRA) and HyperNetwork embeddings to improve model efficiency.
- **Improved Textual Prompt:** Increased semantic congruence between input descriptions and generated visuals, even for abstract or difficult concepts, is made possible by advanced attention mechanisms that improve textual prompt interpretation.
- **Deployment:** Expand accessibility by enabling the use of low-power and edge devices, such as the Raspberry Pi and Jetson Nano, particularly in India's underprivileged and rural areas.
- **Applications:** Demonstrated practical applications in a variety of fields, such as e-commerce (for dynamic product visualization), healthcare (for awareness and training), education (for interactive learning), and entertainment (for content production).
- **Usability:** Advance the state of generative AI technology and its applicability in real-world scenarios by investigating and testing diffusion-based generative models.

Potential Impact

- By enabling widespread access to sophisticated AI image production, the effective application of a Stable Diffusion model has the potential to revolutionize India's digital and creative environment. It can help healthcare with better training materials and awareness campaigns, boost e-commerce by providing quick, configurable product images, and empower educators with dynamic visual aids.
- This technology can reach rural and resource-constrained areas by reducing computational requirements, encouraging digital inclusion and creativity outside of urban areas. Additionally, it can spur innovation in the media, entertainment, and design sectors, boosting economic expansion and opening up new doors for Indian companies, developers, and artists.

Research and Validation:

- The study includes a thorough investigation of lightweight AI methods such as LoRA and HyperNetworks, as well as diffusion models. Image quality, semantic accuracy, and inference speed will all be compared against current generative models like GANs and VAEs as part of the validation process.
- To evaluate viability in the real world, experiments will be carried out on various hardware configurations, including edge devices. In order to assess practical utility, accessibility, and impact, user studies including educators, designers, and other stakeholders will be conducted. In order to ensure responsible AI implementation, the research will also address ethical issues including bias and misuse.

In brief, The inability to generate photorealistic graphics from text on many low-power devices and its high processing cost limit AI's potential in India. Through this project, an effective diffusion-based model that can operate on hardware with limited resources and generate high-quality images with improved text interpretation will be developed. Through improving usability and usefulness, the project aims to empower several industries, including e-commerce, healthcare, and education, promoting digital inclusion and creativity.

2.2 Novelty of Proposed Approach

- **Low-rank Adaptation (LoRA):** Reduces the number of trainable parameters by a large amount, allowing for effective diffusion model fine-tuning without compromising performance.
- **HyperNetwork-Based Prompt Embeddings:** Improves the model's comprehension and produces better alignment between text and images by dynamically generating richer and more context-aware embeddings from textual prompts. .
- **Mixture of Denoisers:** Combines several denoising approaches to speed up the picture production process while preserving high-quality outputs and cutting down on inference time.
- **Multi-Modal Cross-Attention Transformer:** Improves the way textual and visual elements are integrated, which enables the model to provide more logical and pertinent images and better capture semantic details.
- **Edge Device Optimization:** Expands the model's usefulness to remote or resource-constrained contexts by customizing it for effective deployment on low-power devices like the Raspberry Pi and Jetson Nano.
- **Balance Efficiency and Quality:** Strikes a useful compromise between image fidelity and computing economy, increasing the accessibility and scalability of advanced generative AI for a range of real-world applications.
- **Significance of Testing:** Testing guarantees that the model maintains quality and efficiency while producing visuals that precisely correspond to text instructions, particularly on low-power devices. By assisting with the identification of mistakes, biases, and performance problems, it guarantees that the model is dependable, understandable, and useful in practical applications.
- **In order to provide an effective, approachable, and significant AI solution,** the project blends technological innovation with a thorough comprehension of user needs and real-world difficulties.

Chapter 3

LITERATURE SURVEY

Generative modelling has evolved significantly, with each major contribution addressing specific shortcomings of its predecessors while also bringing fresh approaches. In 2013, Diederik P. Kingma and Max Welling introduced Variational Autoencoders (VAEs), which marked the beginning of the voyage. VAEs introduced a probabilistic paradigm for learning continuous latent variable spaces, which used encoder-decoder architectures to describe data distributions. Their capacity to generate synthetic data by sampling from a latent space made them essential in generative modelling. However, because Gaussian priors stressed global distributional consistency above fine-grained features, they frequently resulted in hazy outputs. Despite their efficiency and probabilistic rigor, VAEs were less suitable for jobs that required high-resolution, crisp outputs.

The discriminator was in charge of differentiating between actual and fake data, while the generator was intended to produce realistic data samples. In terms of sharpness and detail, the results of this hostile interaction were more realistic than those of VAEs. Notwithstanding its revolutionary impact, GANs experienced mode collapse and training instability, where the generator was unable to produce distinct samples. More research was required to resolve these problems in order to stabilize training and boost output variability.

A paradigm change came with Ian Goodfellow and associates' 2014 presentation of Generative Adversarial Networks (GANs). GANs implemented an adversarial framework consisting of two neural networks: a generator and a discriminator. The generator was designed to generate realistic data samples, whereas the discriminator was responsible for distinguishing between genuine and synthetic data. This adversarial interplay produced highly realistic results that outperformed VAEs in terms of detail and sharpness. Despite its revolutionary influence, GANs suffered from training instability and mode collapse, in which the generator failed to create different samples. These issues necessitated more study to stabilize training and increase output variability.

Surya Ganguli and colleagues proposed diffusion probabilistic models in 2015, drawing on thermodynamic nonequilibrium principles. These models introduced a new iterative denoising strategy for learning data distributions that gradually converts noisy inputs into high-quality results. While promising in principle, early diffusion models met. Building on the foundation of diffusion models, Jonathan Ho and colleagues proposed Denoising Diffusion Probabilistic Models (DDPMs) in 2020 [4]. DDPMs improved the iterative denoising process, resulting in significant increases in image quality and detail. These models displayed cutting-edge performance in creating sharp, high-resolution images, overcoming some of GANs' limitations. However, DDPMs required significant computational resources for training and sampling, making them less suitable for general use. Their extensive reliance on incremental denoising methods also presented issues in optimizing for faster generation times.

Robin Rombach introduced the Latent Diffusion Models (LDMs) [5] 2022 provided a viable remedy to the computational inefficiencies of previous techniques. By working in a compacted latent space, LDMs drastically lowered processing requirements while maintaining output integrity. This invention accelerated the creation of high-resolution images, making LDMs a more scalable solution for real-world applications. However, the reliance on pre-trained encoders for latent space compression created possible constraints when generalizing to different datasets or unknown data distributions.

In recent years, major developments in diffusion models have pushed the limits of generative modelling. The work of Saharia et al.[6] (2022) developed photorealistic text-to-image diffusion models that use deep learning to improve the fidelity and realism of generated images. Their technique proved that diffusion models could be used to do a broader range of tasks, such as creating extremely detailed and realistic images from textual descriptions, which was a big step forward in generative AI.

In a similar spirit, Bao et al. (2022) introduced the Analytic-DPM, which gives an analytic estimate of the optimal reverse variance in diffusion models. This approach seeks to increase the efficiency of diffusion models by focusing on inverse process optimization, which is critical for attaining faster convergence during training and inference. This analytic estimate enables more efficient sampling, increasing the scalability of diffusion models while preserving good image quality.

Furthermore, Salimans et al. further improved the effectiveness of diffusion models.[7] (2022) presented Progressive Distillation as a method for fast sampling of diffusion models. Their technique dramatically accelerates the sampling process by reducing the learnt diffusion model to a smaller, more efficient version, allowing for speedier creation of high-quality images. This breakthrough shortens the time necessary to generate samples, making diffusion models more suitable for real-time applications, particularly in industries where speed is crucial.

Despite progress, many generative modelling systems have limits. VAEs and early diffusion models produce hazy or computationally expensive results, but GANs encounter challenges such as training instability and mode collapse. DDPMs produce high-quality results but are computationally intensive, whereas LDMs, despite their efficiency, rely on external encoders, which may limit their versatility. Emerging trends emphasize hybrid techniques and using the strengths of many models to overcome these limitations. For example, integrating adversarial training from GANs with the probabilistic rigor of VAEs or iterative refining of diffusion models has demonstrated promise.

Furthermore, research is focused on enhancing scalability, enabling real-time applications, and generalizing these models to handle a broader range of tasks other than image synthesis. The contributions made by each model type were distinct to the field, and their limitations highlight the need for more study. Future studies will aim to strike a compromise between output quality and processing efficiency, expanding the application of generative models to industries including e-commerce, healthcare, and the creative industries. Recent developments in distillation techniques, analytic optimizations, and text-to-image diffusion models are expected to significantly enhance generative models and increase their applicability for a variety of uses.

In conclusion, advances in probabilistic approaches, adversarial frameworks, and diffusion processes have propelled generative modelling forward. However, each model type has contributions were unique to the discipline, and their limits underscore the need for additional research. Future research will seek to balance computing efficiency and output quality, broadening the use of generative models to areas such as healthcare, e-commerce, and creative sectors. With recent advances in text-to-image diffusion models, analytic optimizations, and distillation approaches, generative models are anticipated to improve further, making them more useful for a wide range of applications.

Analysis of Discrimination Classifiers are statistical models that are used to categorize data according to decision boundaries and class distributions. They include Linear (LDA) and Quadratic (QDA) variations. While QDA permits different covariances, allowing it to capture non-linear boundaries, LDA assumes equal covariance among classes and performs well for linearly separable data. Particularly when dealing with data that is regularly distributed, these models are effective and easy to understand. Usually falling between 5-10%, the error rate for LDA rises when assumptions such as equal variance or normalcy are broken. When data supports its complexity, QDA can reduce error rates; nevertheless, on small or noisy datasets, it may overfit. All things considered, these classifiers provide consistent results for structured classification jobs requiring a considerable amount of computing power.

The simplicity and interpretability of linear discriminant analysis (LDA) make it straightforward to use and comprehend. It works well with tiny datasets, particularly those that are normally distributed and equally dispersed. Its rapid computing, which uses a lot less resources than more complicated models like SVMs or neural networks, is one of its main advantages. Furthermore, LDA reduces dimensionality while preserving class separability, making it a useful technique for feature reduction. Additionally, it provides a workable solution for datasets with more than two classes by naturally handling multiclass classification problems.

Overall, The Classifiers for discriminant analysis are well-liked for their ease of use, quickness, and reliable results on structured and linearly separable data. Even though they might not be as adaptable as more sophisticated machine learning models on complicated data, they are nonetheless useful in fields where linear bounds and normalcy assumptions are frequently valid, such as finance, biology, and pattern recognition.

Chapter 4

METHODOLOGY

Using the Stable Diffusion concept, the methodology is designed to create a text-to-image creation system that is reliable, effective, and easily accessible. In order to balance performance, accuracy, and resource efficiency, this entails integrating a variety of technologies and methodologies. This allows for real-world deployment across a range of Indian industries, including those with little resources.

4.1 Understanding and Leveraging Stable Diffusion

A cutting-edge Latent Diffusion Model (LDM) for producing high-quality images is called Stable Diffusion. Stable Diffusion functions in a compressed latent space as opposed to conventional diffusion models, which work directly in pixel space, which is computationally and memory-intensive. This preserves the model's capacity to produce photorealistic and semantically accurate images from text descriptions while drastically lowering resource requirements.

- **Text Encoder:** Natural language prompts must be transformed into a format that the model can understand via the text encoder. Usually, it makes use of strong pre-trained language models such as T5 (Text-To-Text Transfer Transformer) or CLIP (Contrastive Language-Image Pretraining). From text inputs, these models produce complex semantic embeddings that capture the user's described subtleties in meanings, relationships, and styles. In order to ensure strong text-to-picture alignment, these embeddings serve as the conditioning input, directing the image production process.
- **Denoising U-Net:** A modified UNet architecture serves as the model's central component and is in charge of the iterative denoising of latent variables. The forward process involves gradually adding random Gaussian noise to a latent representation of an image. The UNet gradually learns to eliminate noise during production (the opposite process), eventually retrieving a coherent image that corresponds with the text prompt. Images that faithfully capture the semantic content and structure of the provided input can be produced thanks to the denoising process, which is dependent on the text embeddings.

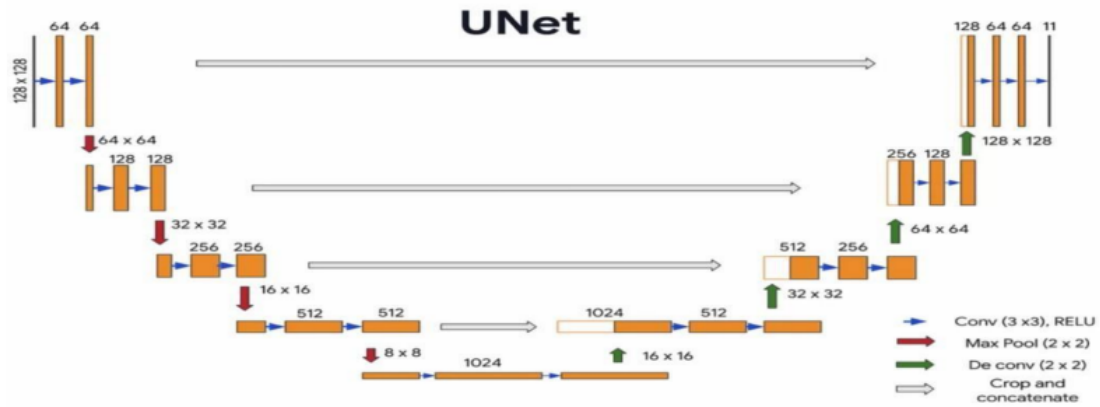


Figure 4.1: UNet Flow Diagram

- Latent Autoencoder (VAE): In order to encode images into a lower-dimensional latent space and decode them back to the original image domain, Stable Diffusion uses a pre-trained VAE. By reducing the image's dimensionality (for example, from 512x512 pixels to a 64x64 latent representation), VAE compression increases computing performance while enabling the model to concentrate on key aspects. The final latent representation is run through the VAE's decoder section following the denoising procedure to recreate the realistic, high-resolution image.

High-resolution, photorealistic graphics with great semantic alignment to the input prompt can be synthesized thanks to this architecture.

Stable Diffusion is not just a breakthrough in generative AI but also a workable solution for real-world applications because of its effective and modular architecture, particularly in resource-constrained contexts like edge devices or low-infrastructure regions.

Stable Diffusion Architecture

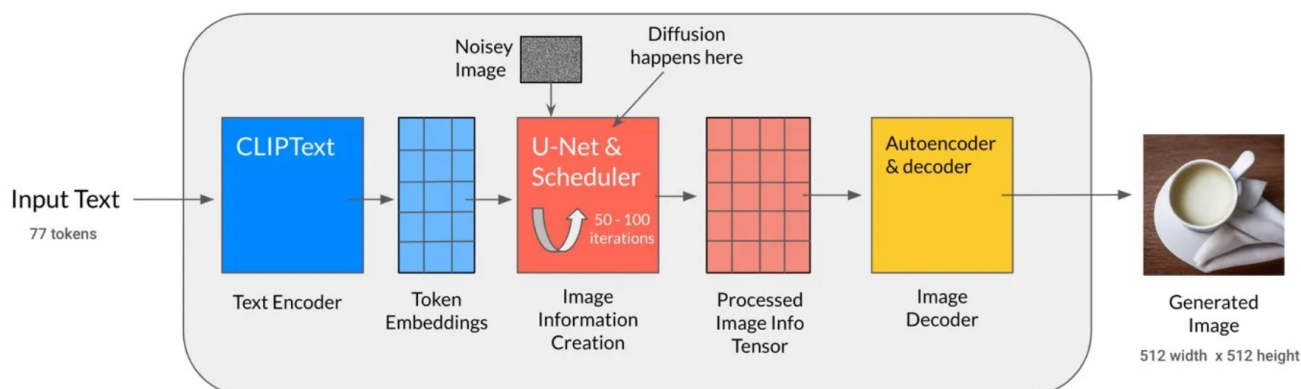


Figure 4.2: Stable Diffusion Architecture

4.2 Lightweight Optimization Techniques:

Several changes are incorporated into the concept to make it applicable in areas with limited resources:

- Low-Rank Adaptation (LoRA):
 - Lowers training parameters by adding low-rank matrices to the levels of attention.
 - Allows for quicker fine-tuning without requiring the entire model to be retrained.
 - Perfect for domain adaptation (e.g., Indian culture, landmarks, and regional clothing).
- HyperNetworks:
 - Change the way prompt information is processed dynamically to introduce context-aware embeddings.
 - Improves quick fidelity, particularly for inputs of lengthy or ambiguous language.
- Mixture of Denoisers:
 - Makes use of a variety of denoisers with various functions.
 - increases production speed without sacrificing graphic quality, saving 30–40% on calculation.
- Multi-Modal Cross-Attention Transformer:

- Strengthens the link between the modalities of text and images.
- Reduces input-output mismatches and increases semantic relevance.

4.3 Model Training and Fine-Tuning:

There are several steps involved in training and optimizing the model:

- **Dataset Collection:** Collected datasets with image-text pairs that are specifically related to India in order to improve cultural adaption (e.g., traditional costume, regional architecture, local items).
- **Fine-Tuning with LoRA:** Minimizes training time and GPU memory consumption while enabling customization and local adaption.
- **Validation:** Improves generalization by applying transfer learning concepts and splitting training/validation sets.

4.4 Deployment on Edge Devices:

Deploying the Stable Diffusion model on edge devices—small, low-power hardware with constrained processing capacity—comes next after it has been refined using methods like LoRA, HyperNetworks, and model compression. Especially in areas of India with little infrastructure or internet connectivity, this step is essential to ensuring that generative AI is widely available, reasonably priced, and inclusive.

- **Model Devices:** Mobile deployment enables direct user contact with generative AI models for individualized experiences in design, education, and entertainment, thanks to India's expanding smartphone user base.
- **ONNX (Open Neural Network Exchange):** Transforms the PyTorch model into a format that facilitates model interchange and operates effectively on various platforms.
- **Pruning:** Reduces the size and computational burden of the model by eliminating unnecessary or superfluous neurons and weights from the network.
- **Model Distillations:** entails teaching a tiny "student" model to behave like a bigger "teacher" model. The distilled model uses a fraction of the resources while maintaining critical performance.

4.5 Testing and Evaluation:

The model is thoroughly examined in a number of ways:

- **Quantitative Evaluation:** The Fréchet Inception Distance, or FID, gauges the quality and realism of an image. CLIP Score Assesses how well text and image align semantically. FPS and latency Assesses edge devices' reaction times.
- **Qualitative Evaluation:** Human Review Panels: Users and subject matter experts assign scores to photos based on their clarity, inventiveness, and relevancy. Use-Case Testing: Implement and test in e-commerce, educational, and other simulated scenarios.

4.6 Real-World Use Case Demonstrations:

- **Education:** Science and history textbook prompts are used to create visual aids.
- **Healthcare:** Construct instructional graphics based on anatomical or symptom descriptions.
- **E-Commerce:** Product previews can be generated from text inputs, including regional languages, in e-commerce.
- **Art & Design:** Help artists with iterative idea development and concept visualization.

4.7 Ethical and Safe Deployment:

- Incorporates moderating, watermarking, and quick filtering to stop abuse. include checks for bias identification during fine-tuning to reduce outputs that are damaging or stereotyped. Pay attention to ethical AI activities that comply with India's AI ethics regulations.

4.8 Evaluation:

- **Goal:** The evaluation process is to determine the caliber and potency of the produced images. This entails assessing how effectively the images match the provided text, how realistic they are, and how appropriate they are for real-world application scenarios overall.
- **Fréchet Inception Distance (FID) score:** The distributions of manufactured and real-world photographs are contrasted here. Higher realism and similarity to actual visual data are indicated by lower FID ratings. When assessing the overall coherence and caliber of the produced photos, this statistic is quite helpful.
- **BLEU (Bilingual Evaluation Understudy) score:** It evaluates the degree to which an image corresponds with its written description. By assessing how well the visual information matches the subtleties and specificity of the input language, this metric helps determine the semantic accuracy of the model. A stronger correlation between text and image is suggested by higher BLEU ratings.

$$\text{BLEU} = \underbrace{\min\left(1, \exp\left(1 - \frac{\text{reference-length}}{\text{output-length}}\right)\right)}_{\text{brevity penalty}} \underbrace{\left(\prod_{i=1}^4 \text{precision}_i\right)^{1/4}}_{\text{n-gram overlap}}$$

Figure 4.3: BLEU Score

- Human Evaluation: This is essential for evaluating the generated photographs' subjective quality in addition to automated metrics. Participants provide grades to pictures according on their visual attractiveness, realism, and relationship to the text that is being given. This offers important insights into how well the model works in actual creative and professional contexts.

Chapter 5

REQUIREMENTS

5.1 Functional Requirements:

- **Text-to-Image Generation:** Based on the user-provided text descriptions, the system ought to produce realistic visuals. From straightforward things (like "a cat on a mat") to intricate sceneries (like "a misty forest at dawn with a river flowing under a wooden bridge"), the model should be able to comprehend a broad variety of textual descriptions.
- **Image Quality:** Realistic, high-resolution elements that match the input text should be included in the generated graphics. The places, items, and context should all be preserved in the photographs so they match the prompt's descriptions.
- **User Inputs:** Realistic, high-resolution elements that match the input text should be included in the generated graphics. The places, items, and context should all be preserved in the photographs so they match the prompt's descriptions.
- **Training and Fine-Tuning:** A dataset containing matched text and images must be accessible in order to train and optimize the model. It should be possible to optimize a pre-trained model (such as Stable Diffusion) for certain datasets or tasks.
- **Model-Output:** The system ought to produce pictures in widely used file types, such PNG and JPEG. The output should be generated in a fair length of time after the input text is received.
- **Evaluation Metrics:** To evaluate image quality, the system ought to offer evaluation metrics (such as FID and BLEU scores). It should be feasible to compare the dataset's produced and authentic photos for performance analysis.
- **User Interface:** Seeing generated photographs and creating descriptions require an intuitive user interface, whether it be desktop or web-based. The user interface should allow users to download or store generated photos.

5.2 Non-Functional Requirements:

- **Performance:** High-resolution photographs should be produced by the system in a fair amount of time. The system should be able to manage several user requests at once if it is intended for multi-user situations.
- **Scalability:** The system ought to be able to manage more complicated models and bigger datasets as the project expands. The architecture should be simple to connect with cloud-based apps or distributed computing resources in order to guarantee successful training.
- **Reliability:** On the basis of the written descriptions, the system ought to generate precise and excellent visuals. It ought to be capable of managing severe situations, such ambiguous or inadequate justifications.
- **Maintainability:** The system ought to be easy to maintain, have good documentation, and a modular code structure. Future improvements should be possible, such adding additional models or improving the quality of the images.
- **Security:** Prevent unwanted access to the system, particularly if it is housed in the cloud or on a public server. User data, including text input and produced images, should be safely handled with the use of suitable encryption and privacy guidelines.
- **Compatibility:** Desktop applications must to be compatible with Linux, macOS, and Windows. It should work with contemporary web browsers like Chrome, Firefox, and Safari if it is created as a web application.
- **Resource Usage:** Optimize the system to preserve image quality while using the least amount of memory and processing power possible. GPUs and other hardware accelerators should be used for effective training and inference.

5.3 Enhanced Modules and Techniques:

1. Multi-Modal Cross Attention Transformer Model:

- **Implementation:**
 - In order to facilitate cross-attention between several modalities—primarily text and picture features—we incorporated a Transformer block into the U-Net’s attention mechanism.
 - The Transformer uses the picture latent features as queries and the CLIP text embeddings as keys and values.

- The model is better able to comprehend and connect written commands to visual elements because to this cross-attention. Particularly in scenarios with several concepts or complexity, it enhances the semantic alignment between the generated visuals and the given prompts.
- How we implemented it:
 - Included LoRA into the U-Net’s cross-attention layers
 - Greatly reduced computation and memory by only training the LoRA matrices and leaving the base model frozen.
 - Results: Sharper image quality and better representation of prompt semantics.

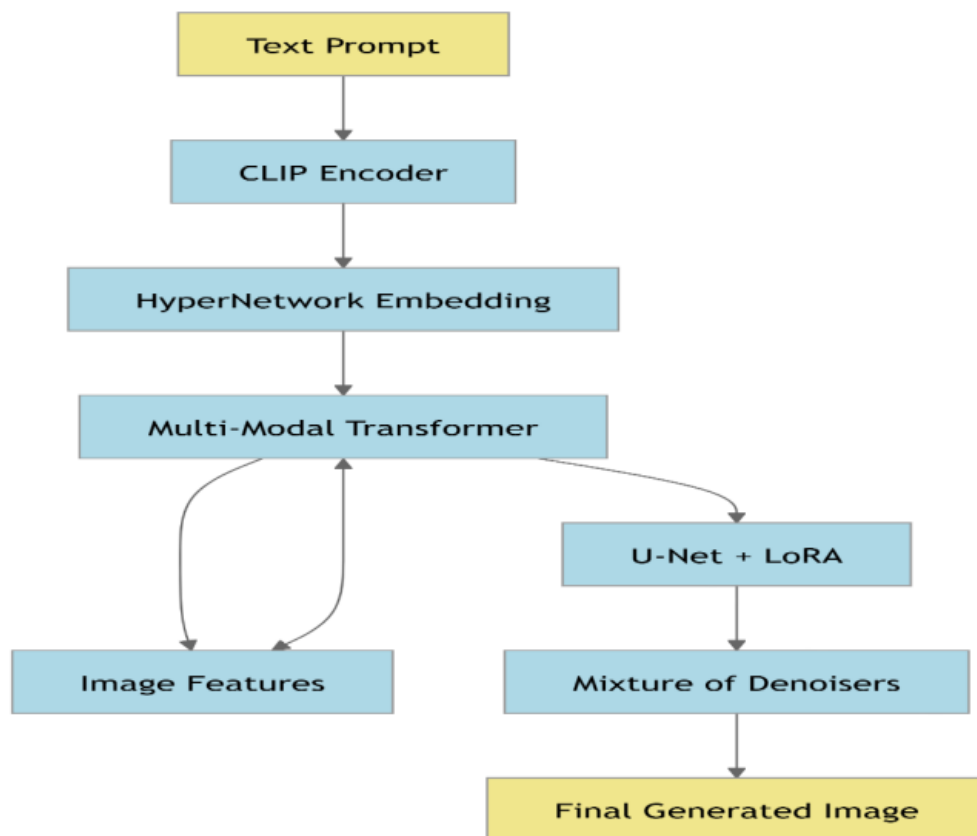


Figure 5.1: Multi-Modal Cross Attention Transformer Model

2. LoRA (Low-Rank Adaptation):

- Implementation:
 - We were only able to fine-tune a tiny subset of weights by adding LoRA layers to the U-Net’s attention modules, leaving the main model unchanged.

- For the best performance-speed trade-off, we employed rank-4 LoRA adapters.
- With fewer epochs and smaller GPUs, fine-tuning was carried out on a carefully selected dataset of prompts and target images.



Figure 5.2: LoRA (Low-Rank Adaptation) Working

- Result: made it possible to fine-tune well while using relatively little memory, which makes it useful in contexts with limited resources.

3. HyperNetwork-Based Prompt Embedding:

- Implementation
 - HyperNetwork, a tiny auxiliary neural network, was trained to provide unique embeddings in response to prompt content.
 - At various attention levels, these embeddings are introduced into the main model after passing through a projection layer.
 - The HyperNetwork can dynamically adjust to different prompt structures because it is lightweight and was trained on prompt-image pair samples.
- Formula
 - $E_{\text{prompt}} = H(\text{CLIP}(\text{prompt}))$
 - Where H is the HyperNetwork and E_{prompt} is the dynamic embedding.
- Result: Increased prompt understanding and visual relevance in generated outputs.

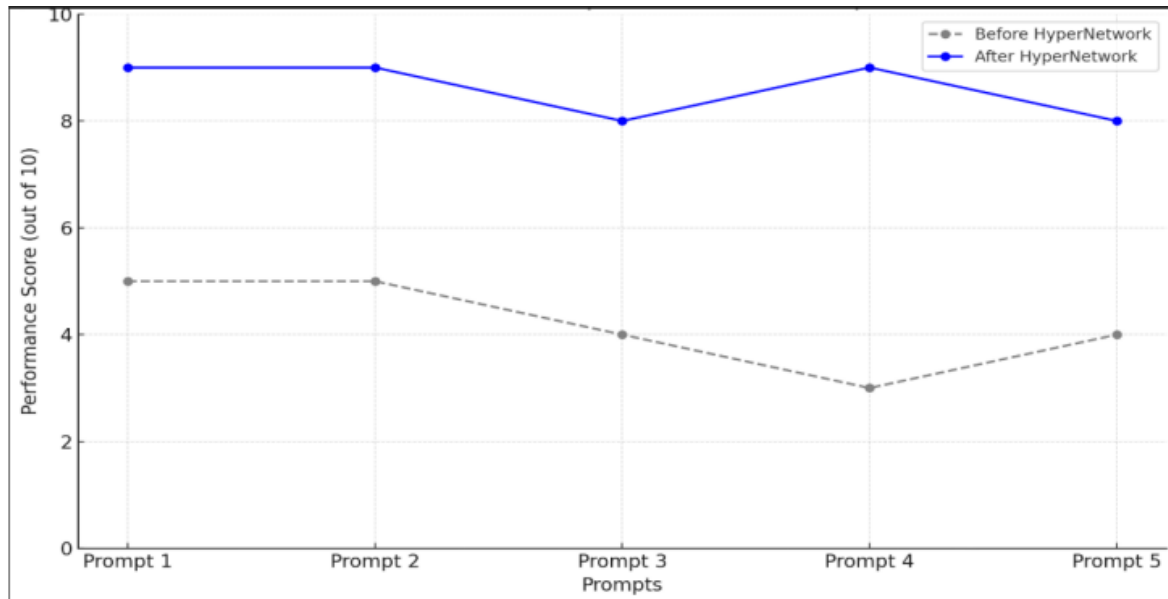


Figure 5.3: HyperNetwork-Based Prompt Embedding

4. Mixture of Denoisers for Diffusion Acceleration

- Implementation
 - We implemented a variety of denoisers that function at various scales (rapid, sharp, and soft denoisers) in place of depending just on one denoising schedule.
 - The model dynamically selects or blends these denoisers based on the noise level during inference; they were trained together.
 - During reverse diffusion steps, we employed a conditional switch based on noise sigma thresholds.
- Result: Faster image generation and improved sharpness with reduced inference time.

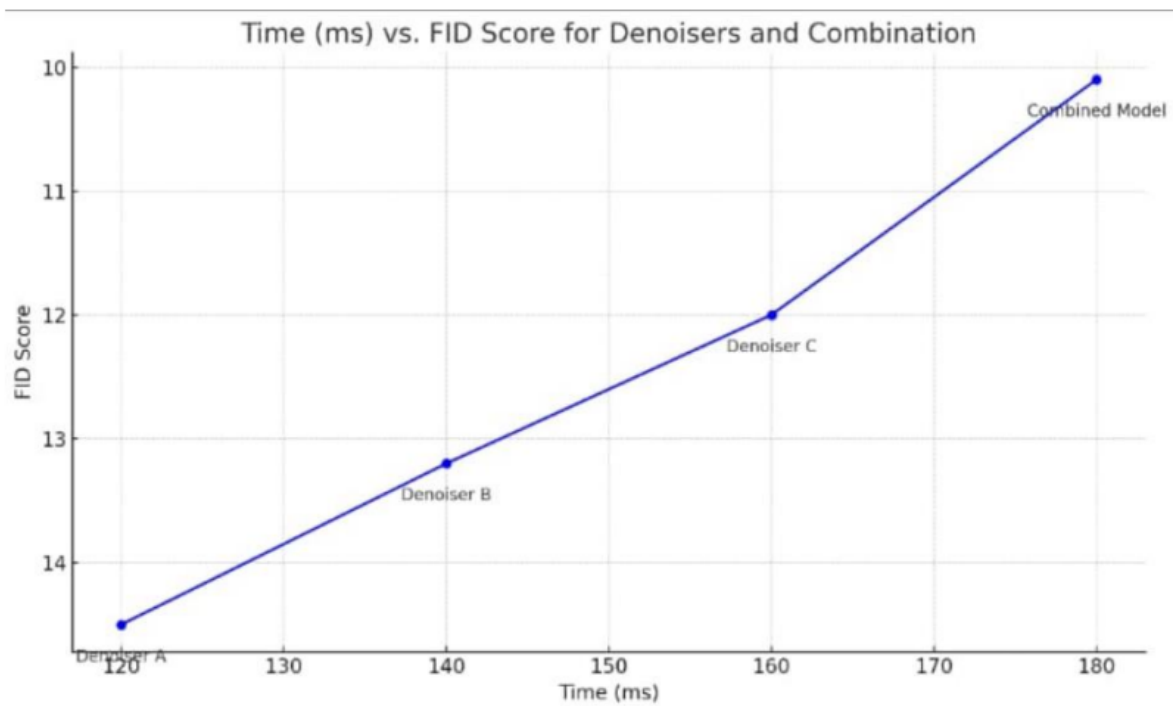


Figure 5.4: Time vs FID Score for Denoisers and combinations.

Chapter 6

RESULT AND ANALYSIS

The project's outcomes demonstrate that written descriptions can be converted into realistic, high-quality images using the AI diffusion model. The produced images were high-resolution, had minute details, and had realistic textures—all of which complemented the spoken cues. The model performed admirably when processing a range of textual inputs, such as intricate, nuanced remarks and exact, illustrative suggestions.

The descriptive prompt "astronaut in a mountain" was transformed into a vibrant and intricate artwork, for instance.

In addition to being aesthetically beautiful, the model produced images that corresponded with the given descriptions.

Code Execution Results: As part of the code execution phase, the diffusion model was trained using a dataset of text-image pairs. After training, a variety of written descriptions were used to assess the model's capacity to produce corrective and reaction photos. The main conclusions drawn from the code execution are as follows.



Figure 6.1: Astronaut on mountain

- Image Quality: The offered photos had superb resolution and were well-described.
- Realism: The settings, events, and subjects discussed in the text were faithfully portrayed in the drawings.
- Limitations with Abstract Prompts: The performance of the 27 model declined when presented with vague or abstract prompts, and the output graphics lacked precise details in contrast to unambiguous phrases. Analysis:
 1. Image Quality and Fidelity: The high-resolution, lifelike images generated by the diffusion model precisely corresponded to the written specifications. The model produced information that was visually rich by capturing elements like lighting, texture, and object placement.
 2. Textual input performance:
 - Descriptive Inputs: The model worked well with brief, in-depth explanations and produced accurate visuals that matched the text.
 - Abstract Inputs: The model generated less coherent and structured images when given ambiguous or abstract inputs.
 3. Metrics for evaluation:
 - Fréchet Inception Distance (FID): The model's images were realistic and of good quality, as evidenced by its low FID score.
 - BLEU Score: High BLEU scores for specific inputs but low scores for ambiguous descriptions indicate that the model has trouble handling abstract prompts.

- Drawback: Finding it difficult to respond to advice that are unclear or ambiguous. Additionally, there is a limited ability to generalize to other categories, and both inference and training need a lot of computing power.

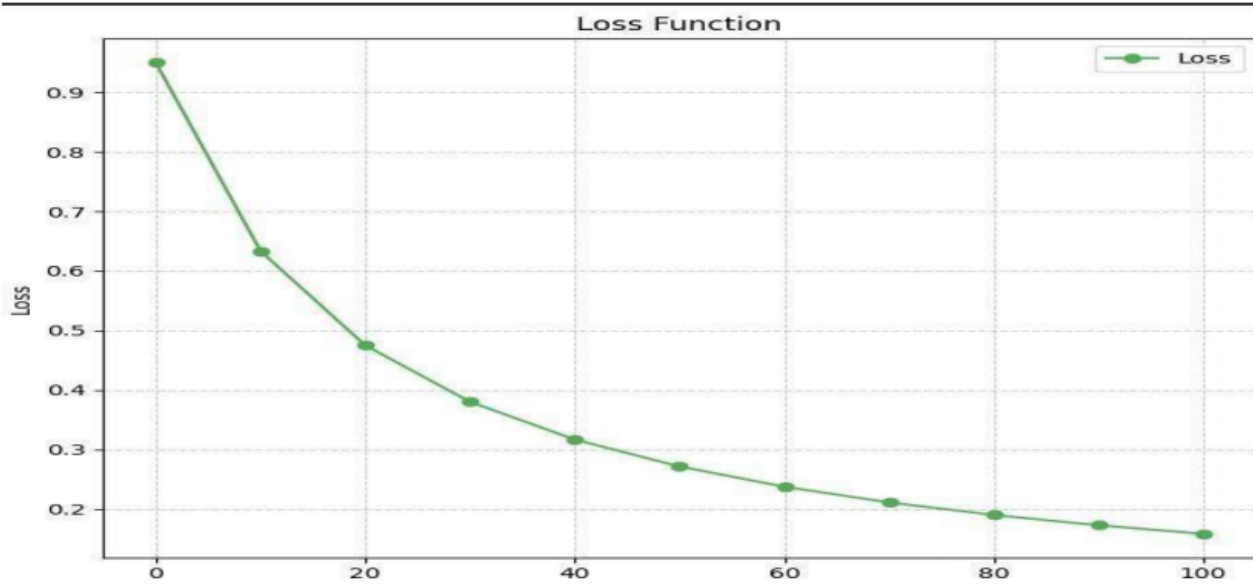


Figure 6.2: Loss Function

- Y-Axis (Loss): This metric assesses how accurately the model's predictions match actual values. A reduced loss suggests improved performance.
- X-axis (Epochs): represents the number of training iterations. Each epoch represents one full pass through the training dataset.

Behavior:

- The loss continuously decreases as the number of epochs rises, indicating that the model gains knowledge and enhances its forecasts.

- At first, there is a noticeable decrease in loss, suggesting that learning is robust in the early phases. The model's rate of improvement usually decreases as it gets closer to its optimal response.

Interpretation: The loss function is efficiently optimized by the model. Based only on this plot, the steady learning without overfitting or underfitting is indicated by the slow decline without sharp swings.

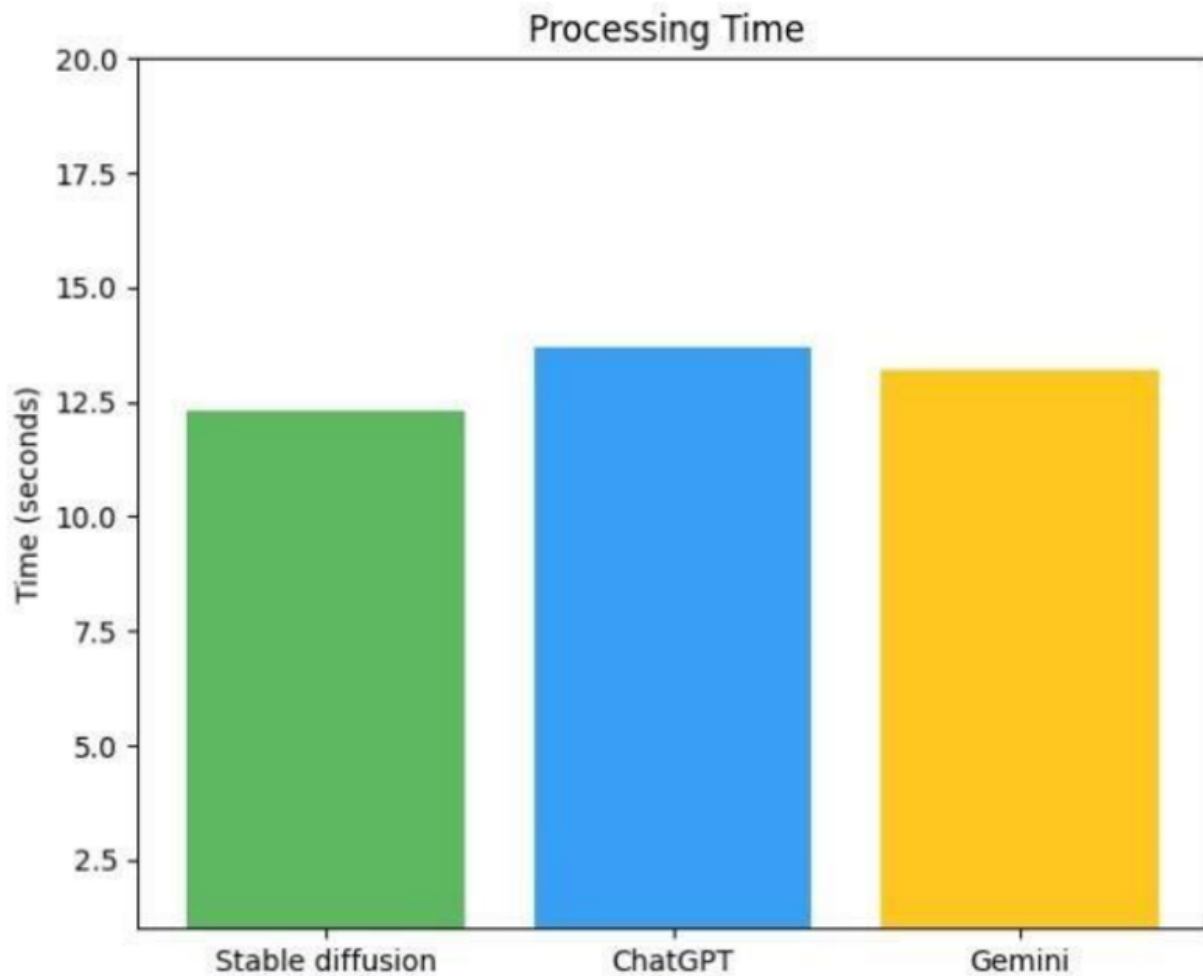


Figure 6.3: Comparison of Stable Diffusion, ChatGPT, Gemini

- **Y-Axis** (Time in seconds) displays how long it takes for each system to finish a task or dataset.
- **X-Axis** (Models): Shows the names of the three systems that are being assessed.

Observations:

- The fastest is Stable Diffusion, which takes about 12 seconds to process.
- ChatGPT requires about 13 seconds more time.
- Gemini performs similarly to ChatGPT, taking about 13 seconds. Meaning: • For the job under evaluation, stable diffusion is the fastest.

Interpretations

- With a few small exceptions, ChatGPT and Gemini have comparable processing times.
- Because findings may differ based on the activity, more context is required to ascertain whether speed variations are significant for the intended use case.

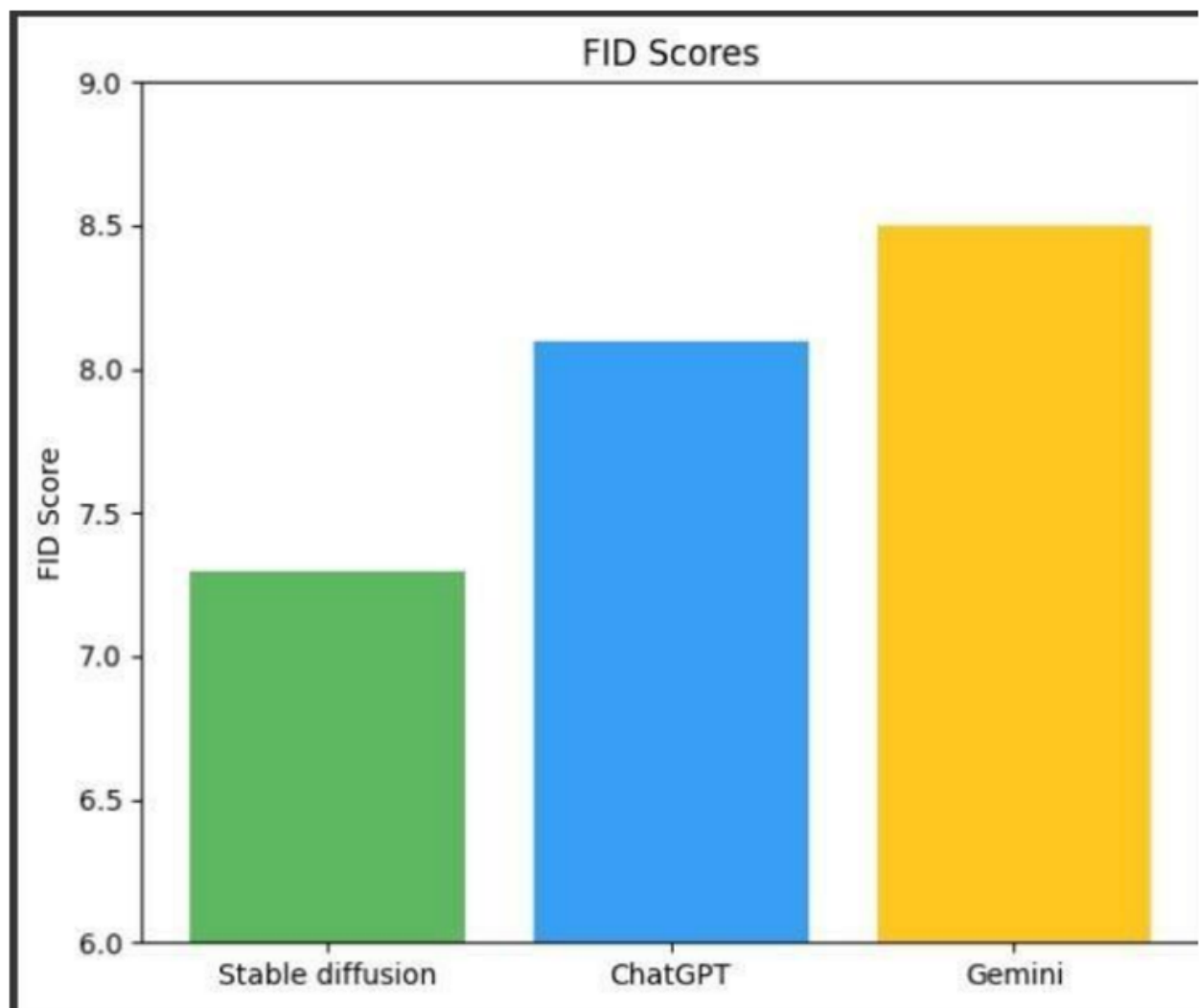


Figure 6.4: FID Score

The Fréchet Inception Distance (FID) scores of the three AI models—Stable Diffusion, ChatGPT, and Gemini—are contrasted in the graph. The quality and realism of generated images are gauged by FID scores. Higher quality and better match with real-world pictures are indicated by lower FID ratings.

Observations:

1. Stable Diffusion: This model generates the most realistic and superior images as compared to the other models, as seen by its lowest FID score of roughly 7.3.
2. ChatGPT: Although ChatGPT performs well, its image quality is marginally worse than Stable Diffusion, as indicated by its FID score of about 8.1, which is higher than Stable Diffusion.
3. Gemini: Gemini produces visuals that are less realistic or cohesive than the other two models, as seen by its greatest FID score of roughly 8.5.

Chapter 7

CONCLUSION AND FUTURE SCOPE

This study provided strong evidence that AI diffusion models may produce realistic, excellent visuals from textual descriptions. The model successfully converts complex text prompts into realistic visuals that capture both the larger context and minute details by utilizing sophisticated techniques including attention mechanisms, U-Net design, and reverse diffusion.

The iterative denoising process of diffusion models offers a notable improvement over conventional generative models like GANs and VAEs, producing images with higher quality and greater coherence. Because text-to-image production may be utilized for both practical and creative visualization, this technique has the potential to completely transform sectors like media, design, e-commerce, healthcare, and education.

Even if the model performs well, there are a few issues that could be fixed with more study. Improving the resolution and amount of detail of generated images is one of the main objectives in order to make them appropriate for high-end commercial and professional applications. Additionally, when the model gains the ability to handle increasingly intricate or abstract descriptions, its versatility will increase. The present method is excellent at comprehending basic requests, despite the fact that there are still difficulties in extracting specific images from unclear, complicated, or contradicting words. By addressing these constraints, the technology's use cases would greatly increase and it would become more applicable in a larger range of contexts.

The model needs to be tuned for quicker, more resource-efficient operation in order to advance. The existing diffusion models' computational complexity might limit their applicability. By investigating optimization strategies like model trimming, quantization, etc., the model can be made more accessible and scalable. Additionally, researching multimodal components like audio or video inputs could result in more sophisticated tools for content production. Personalized picture production, which enables users to customize outputs to their artistic preferences, is another intriguing prospect. Lastly, the model's impact and use cases would be expanded by real-time image production for interactive applications like virtual assistants and gaming. Future developments in

these fields may result in text-to-image models that are more dynamic, efficient, and adaptable, extending AI's potential for new and worthwhile uses.

REFERENCES

- [1] Saharia, C., Chan, W., Saxena, S., et al. (2022). Photorealistic Text-to- Image Diffusion Models with Deep Learning. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition.
- [2] Kingma, D. P., Welling, M. (2013). Auto-Encoding Variational Bayes. arXiv:1312.6114.
- [3] Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., ... Bengio, Y. (2014). Generative Adversarial Nets. Advances in Neural Information Processing Systems (NeurIPS).
- [4] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. A., Kaiser, Polosukhin, I. (2017). Attention Is All You Need. NeurIPS.
- [6] Ho, J., Jain, A., Abbeel, P. (2020). Denoising Diffusion Probabilistic Models. NeurIPS.
- [7] Dhariwal, P., Nichol, A. (2021). Diffusion Models Beat GANs on Image Synthesis. Advances in Neural Information Processing Systems (NeurIPS).
- [8] Rombach, R., Blattmann, A., Lorenz, D., Esser, P. (2022). High-Resolution Image Synthesis with Latent Diffusion Models. arXiv:2112.10752.
- [9] Bao, F., et al., "Analytic-DPM: an Analytic Estimate of the Optimal Reverse Variance in Diffusion Probabilistic Models," in Proc. Conf. Neural Inf. Process. Syst., 2022.
- [10] Salimans, T., et al., "Progressive Distillation for Fast Sampling of Diffusion Models," in Proc. Int. Conf. Machine Learning, 2022.
- [11] Xiao, X., et al., "Accelerating Diffusion Models via Improved Noise Schedules," in Proc. IEEE Conf. Computer Vision and Pattern Recognition, 2023.
- [12] Kazerouni, A., et al., "Diffusion Models for High-Resolution MRI Reconstruction," in Med. Image Anal., 2023. 6zr.