

DECLARATION

The project entitled “**TEXT-TO-IMAGE GENERATOR**” is being submitted to the Department of Computer Science and Information Technology, Dillibazar, Kathmandu, Nepal for the fulfillment of the final year project under the supervision of Mr. Avishek Kuinkel. This project is original and has not been submitted earlier in part or full in this any other form to any university or institute, here or elsewhere, for the award of any degree.

By

Neha Shrestha, 24287/076

Norden Ghising Tamang, 24290/076

RECOMMENDATION

This is to recommend that **Neha Shrestha and Norden Ghising Tamang** have carried out research entitled “**Text-to-Image Generator**” for the fulfillment of the seventh semester in Bachelor’s degree of Computer Science and Information Technology under Tribhuvan University under my supervision. To my knowledge, this work has not been submitted for any other degree.

They have fulfilled all the requirements laid down by the Trinity International College Department of Computer Science and Information Technology, Dillibazar, Kathmandu, Nepal.

Avishek Kuinkel

Project Supervisor,

Department of Computer Science and Information Technology,

Trinity International College

Dillibazar, Kathmandu, Nepal

ACKNOWLEDGEMENT

We would like to express our deepest gratitude to the people who supported us throughout our project. Firstly, we want to thank **Mr. Abhishek Dewan**, Assistant Program Coordinator at Trinity International College, for providing us with guidance, support, and supervision throughout the project. His insights and feedback have been invaluable in shaping this project documentation.

Furthermore, we are grateful to **Mr. Avishek Kuinkel**, who acted as our Project Supervisor, for his suggestions and directions that helped us throughout the project. His wise advice and suggestions enabled us to adapt to the growing obstacles we encountered.

We also want to thank all the instructors in the Department of Computer Science and Information Technology for their unwavering support, encouragement, and direction. Their assistance was crucial to our project's successful advancement.

Finally, we would like to express our gratitude to the friends and department personnel that helped us finish the project successfully. We sincerely appreciate all their assistance.

Neha Shrestha, 24287/076

Norden Ghising Tamang, 24290/076

ABSTRACT

The “Text-to-Image generator” presents an innovative approach to synthesize visual content from textual descriptions. Leveraging the power of LDMs, the system employs a novel iterative process to refine a latent representation space, enabling the generation of high-quality images from diverse textual prompts. The project aims to provide users with a versatile tool for creative expression, educational illustration, and efficient content creation by combining the strengths of latent diffusion models and convolutional neural networks in a unified framework. This project uses advanced Generative AI techniques like the diffusion model, UNet, and Variational Autoencoder to generate visual outputs.

The project's core focus lies in image generation, where textual prompts act as creative guides the generation process. By incorporating the diffusion model, the system undergoes a refined process of noise reduction in a latent space, resulting in the production of images. Variational Auto-encoders used for latent image manipulation and UNet is used for prediction of noise. Through iterative Markov process, an image is generated.

The implementation employs the PyTorch framework along with the Diffusers library ensuring a flexible and efficient development environment. These frameworks provide specialized toolkit for managing the diffusion process, optimizing performance, and enhancing the overall stability of the image generation pipeline.

Keywords: *UNet, Variational Auto-encoders, DDPM, CNN, Python Programming, PyTorch, Diffusers*

TABLE OF CONTENTS

DECLARATION	ii
RECOMMENDATION	iii
LETTER OF APPROVAL	iv
ACKNOWLEDGEMENT	v
ABSTRACT	vi
LIST OF FIGURES	ix
LIST OF ABBREVIATIONS	x
CHAPTER – 1: INTRODUCTION	1
1.1 Introduction	1
1.2 Problem Statement	1
1.3 Objective	2
1.4 Scope & Limitations	2
1.5 Development Methodology	2
1.6 Report Organization	3
CHAPTER – 2: BACKGROUND STUDY AND LITERATURE REVIEW	4
2.1. Background Study	4
2.2. Literature Review	4
CHAPTER – 3: SYSTEM ANALYSIS	7
3.1 System Analysis	7
3.1.1 Requirement Analysis	7
3.1.2 Feasibility Analysis	9
3.1.3 Analysis	10
CHAPTER – 4: SYSTEM DESIGN	11
4.1 Design	11
4.1.1 Sequence diagram	11
4.1.2 Activity diagram	12
4.2 Algorithm Details	13
CHAPTER – 5: IMPLEMENTATION AND TESTING	18
5.1 Implementation	18
5.1.1 Tools Used	18
5.1.2 Implementation Details of Modules	18
5.2 Testing	21
5.2.1 Unit Testing:	21
5.2.2 System Testing:	22

5.2.3 User Testing:	23
5.3 Result Analysis	23
5.3.1 Loss Analysis:	23
CHAPTER – 6: CONCLUSION AND FUTURE RECOMMENDATIONS	24
6.1 Conclusion	24
6.2 Future Recommendations	24
REFERENCES	25
APPENDIX-I	26
APPENDIX-II	36

LIST OF FIGURES

Chapter - 3:

Figure 1: Use Case Diagram.....	8
Figure 2: Gantt-Chart of the project.....	9
Figure 3: Flow Diagram of the system.....	10

Chapter - 4:

Figure 4: Sequence Diagram.....	11
Figure 5: Activity Diagram.....	12
Figure 6: Latent Diffusion Architecture [1].....	13
Figure 7: VAE Architecture.....	13
Figure 8: VAE Down Sample, VAE Block, VAE Up Sample (Left-to-right)	14
Figure 9: VAE Down and Up Block	14
Figure 10: Time Embedding and Condition Embedding	14
Figure 11: UNet Architecture [4].....	15
Figure 12: UNet for LDM.....	16
Figure 13: Conv In, Resnet, Conv Out (Left-to-right)	16
Figure 14: UNet: Down Sample, Up Sample, Down Block, Mid Block, Up Block.....	16
Figure 15: Training and Sampling Algorithm [3]	17

Chapter - 5:

Figure 16: VAE Test Case. Left (Image), Right(Latent).....	21
Figure 17: U-Net Noise Prediction Test. Left (Actual), Right (Predicted).	22
Figure 18: System Test – I.....	22
Figure 19: System Test - II	22
Figure 20: Loss Curve (MSE Loss).....	23

LIST OF ABBREVIATIONS

AEVB	Auto-Encoding Variational Bayes
CE	Conditional Embedding
CNN	Convolutional Neural Networks
DDPM	Denoising Diffusion Probabilistic Models
GAN	Generative Adversarial Network
LDM	Latent Diffusion Model
MSE	Mean Squared Error
SGVB	Stochastic Gradient Variational Bayes
SiLU	Sigmoid Linear Units
TE	Time Embedding
VAE	Variational Auto-encoders
VLB	Variational Lower Bound
VQGAN	Vector Quantized Generative Adversarial Network

CHAPTER – 1: INTRODUCTION

1.1 Introduction

Generative AI, an evolving paradigm within artificial intelligence, showcases a creative and innovative part of technology that empowers machines to produce content resembling existing data. This versatile technology spans various mediums including images, text, audio, video and more. The concept behind this mechanism involves extracting and understanding patterns and information present in the data on which the models are trained.

One of the interesting generative AI models is the LDM. In contrast to the previously defined models, this model does not rely on extensive spatial compression. It operates within a learned latent space which is characterized by superior scaling properties concerning spatial dimensionality [1]. In simpler terms, the model doesn't compromise image quality or information while compressing spatial aspects, ensuring a more effective representation of the data. The main idea behind this is to systematically decompose the data structure through an iterative forward diffusion process. Then a reverse diffusion process is applied that restores structure in the data, yielding a highly flexible and tractable model of the data. This diffusion probabilistic (diffusion) model operates as a latent text-to-image diffusion model, repeatedly diminishing noise in a latent representation space and then transforming that representation into a complete image [3].

This project proposes a model that shows the interplay of noise and neural networks to produce compelling results. Thousands of images are meticulously registered with layers of noise, forming a corpus of pure noise images. Subsequently, these noisy images undergo training in a Markov chain process to revert to their original state through neural networks [3]. The machine learns to discern and eliminate this noise so adeptly that the model becomes proficient in transforming any arbitrary noisy input into a new image that aligns with our training data. This process empowers the model to generate the desired image corresponding to any textual prompt.

1.2 Problem Statement

Generative AI emerges as a solution to the challenge of data scarcity by enabling the creation of vast training datasets. Moreover, overcoming communication barriers associated with novel concepts, the text-to-image model serves as a facilitator for improved understanding and learning. Furthermore, this model expedites the transformation of ideas into visual

representations, offering a swift and intuitive process devoid of technical complexities. In essence, generative AI not only addresses data limitations but also acts as a catalyst for enhanced communication and streamlined visual expression, catering to a diverse range of users who may encounter difficulties in grappling with new information.

1.3 Objective

The objectives of this project are:

- To transform textual descriptions into visually appealing images.
- To ensure the generated images are coherent and contextually accurate.
- To provide a user-friendly interface for users to input text and retrieve generated images.

1.4 Scope & Limitations

The scopes of this project are:

- The project can be used to express textual prompts to visual images.
- The system can serve as an educational tool, aiding in the creation of visual aids for educational materials.
- Preferred text prompts can be given by users to derive images.
- The project facilitates easy user interaction.

The limitations of this project are:

- The project's limitation lies in its ability to generate only a limited class of images (2-5).
- The quality and diversity of generated images heavily depend on the training data.
- The system might struggle with highly complex and abstract textual prompts.
- Without a powerful GPU system, processing is slow and scalability is limited.

1.5 Development Methodology

The developmental methodology for this project involves conducting a thorough literature review on LDM and Diffusion Probabilistic Models. Further, the procedures are collecting and preparing relevant datasets, adding noise to the datasets, designing U-Net with proper attention mechanism and embeddings of time and conditions, implementing VAEs, executing the model using different frameworks, evaluating and validating the model's performance, iteratively

refining the model based on the feedback and evaluation metrics, and documenting the entire development process in a comprehensive report.

1.6 Report Organization

The report is organized as follows:

Chapter 1 introduces the concept of Generative AI and its ability to craft diverse content formats. The focus here is on probabilistic models that can generate high-quality images and propose the use of latent diffusion models. This chapter outlines the problem statement, objectives, scope and limitations of the project. Overall, chapter 1 lays the foundation for the rest of the report by briefly addressing the goals of the project.

Chapter 2 of the report contains background information and a literature review on the topic of the diffusion probabilistic model and its utilization in the generation of sophisticated images. The literature review covers various studies on the use of a combination of CNNs and transformers, employing pre-trained autoencoders and applying cross-attention mechanisms. The chapter provides a summary of the major findings and contributions of each study.

Chapter 3 of the report is dedicated to system analysis. This section covers the study of several analyses such as requirement analysis, feasibility analysis and project analysis. The requirement analysis includes both functional and non-functional requirements. The feasibility analysis examines technical, schedule and economic feasibility. Moreover, the analysis section provides a flow diagram that illustrates the project's workflow. Specifically, a flow diagram depicts the step-by-step process of the project, providing a concise understanding of how the project works.

Chapter 4 of this report emphasizes the design aspect of the project. This commences with a sequence diagram that outlines the sequence of interactions between the system and the user. Following that is an activity diagram that exemplifies the various stages that occur within the system. Additionally, this chapter includes a detailed description of the algorithm used in the project.

Chapter 5 describes the implementation and testing part of the project. The tools and algorithms used to develop the application are covered in this chapter, which also focuses on testing the application with various test case types.

The conclusion and future improvements are covered in **Chapter 6**, which offers suggestions on what we accomplished at the project's finish and how we might improve the application going forward.

CHAPTER – 2: BACKGROUND STUDY AND LITERATURE REVIEW

2.1. Background Study

Diffusion models are a class of generative models in machine learning that simulate the data generation process by transforming a simple and easily sample able distribution, typically a Gaussian distribution, into a more complex data distribution of interest. Diffusion models were inspired by the natural diffusion process in physics, which describes how molecules move from high-concentration to low-concentration areas.

The idea of diffusion models was first introduced in a 2015 paper named "The Deep Unsupervised Learning using Nonequilibrium Thermodynamics" and gained momentum in 2020 with the publication of several papers [6]. In this paper, the approach to achieving both flexibility and tractability is to systematically and slowly destroy data structure through an iterative forward diffusion process followed by the reversal of a Markov diffusion chain. The result is an algorithm that can learn a fit to any data distribution which is straightforward to manipulate conditional and posterior distributions.

Diffusion models have diverse applications across several domains, such as text-to-video synthesis, image-to-image translation, image search, and reverse image search. They have better image quality, interpretable latent space, and robustness to overfitting compared to traditional generative models.

2.2. Literature Review

The article "Denoising Diffusion Probabilistic Models" is built on top of the 2015 paper "The Deep Unsupervised Learning using Nonequilibrium Thermodynamics". Here, the authors took inspiration from nonequilibrium thermodynamics and found a correlation between diffusion probabilistic models and denoising score matching with Langevin dynamics. They have presented a class of latent variable models that use a parameterized Markov chain to gradually add noise to the data in the opposite direction of sampling. The models are trained using a forward diffusion process that maps data to noise and a learnt, parametrized reverse process that performs iterative denoising, starting from pure random noise. The parameterization of the diffusion models claims to be the primary contribution to the best sample quality results. They

make a clear statement that their models do not have competitive log-likelihoods compared to other models. The sampling procedure of the diffusion model is discovered to be progressively decoding often resembling autoregressive decoding [3].

A few simple modifications are made on top of the above paper to achieve an improved version of DDPM. In this paper, DDPMs can achieve log-likelihoods with other likelihood-based models even on highly diverse datasets like ImageNet. The authors discovered a learnable variance schedule in the reverse process using simple reparameterization and hybrid learning objectives that combine the VLB with a simplified objective from the above paper. They found that while the linear noise schedule used in the previous paper worked well for high-resolution images, it was sub-optimal for images of resolution 64×64 and 32×32 . Therefore, to address the problem of a forward noising process being too noisy, they have constructed a cosine schedule which has a more gradual loss. Incorporating all the requirements, they realized that pre-trained hybrid models could achieve good samples with as few as 50 forward passes if the previous DDPM required 100s of forward passes to produce good samples. They have additionally used precision and recall to compare the distribution coverage between DDPMs and GANs, finally noting that the former covered a larger portion of the target distribution [5].

The efficiency of synthesizing high-resolution images is maximized by integrating the usage of transformers as well as CNNs. The power of the inductive bias of CNNs with the articulation of transformers enabled to model and thereby produce quality images. The approach is to use a convolutional VQGAN to learn a codebook of context-rich visual parts, whose composition is subsequently modelled with an autoregressive transformer architecture. A discrete codebook provides the interface between these architectures and a patch-based discriminator enables strong compression while retaining high perceptual quality. The study was bounded up to only 16×16 discrete latent spaces which produced only 256 tokens thus, no super big image could be taken as input otherwise reconstruction would have degraded drastically. Therefore, the researchers used the concept of a sliding attention window which solves the lower computation limitation of the project. They have also replaced MSE loss with perceptual loss and adversarial loss [2].

With diffusion models explained as a sequential technique, the typical operation is performed in pixel space. The formulation algorithm executes by adding or removing noise to a tensor of the same size as the original image resulting in slow inference speed and high computational cost.

Thus, the paper “High-Resolution Image Synthesis with Latent Diffusion Models” breaks the ice and deals with issues of previous approaches through the use of latent space of powerful pre-trained autoencoders. The researchers used the models that can be interpreted as an equally weighted sequence of denoising autoencoders which can be trained to predict a denoised variant of their input. The cross-attention conditioning mechanism is used to train a large 1.4 billion parameter text image diffusion model. This model consists of the U-Net and the transformer backbone, jointly trained on the publicly available LAION 400M dataset. The resulting model can compose samples from complex text prompts and write user-specific text [1].

The paper “Auto-Encoding Variational Bayes” addresses the challenge of performing efficient inference and learning in directed probabilistic models with continuous latent variables and large datasets. It introduces a stochastic variational inference and learning algorithm, SGVB, which scales effectively to large datasets and can even handle intractable posterior distributions given certain differentiability conditions. The key contributions lie in two aspects: firstly, a reparameterization of the variational lower bound leads to a lower bound estimator that can be optimized using standard stochastic gradient methods, and secondly, for independent and identically distributed datasets with continuous latent variables, the proposed algorithm, AEVB, efficiently fits an approximate inference model to the intractable posterior using the SGVB estimator. The efficacy of the approach is demonstrated through theoretical analyses and experimental results [7].

The U-Net architecture is introduced, which utilizes data augmentation to maximize the use of annotated samples. It features a contracting path for context capture and an expanding path for precise localization. The network is shown to be highly effective, surpassing previous methods on the biomedical imaging challenge for neuronal structure segmentation and cell tracking. Notably, it demonstrates superior performance and computational efficiency, even across diverse biomedical segmentation tasks like transmitted light microscopy. With minimal annotated data and data augmentation, particularly with elastic deformations, U-Net achieves remarkable results, making it a versatile tool for biomedical segmentation tasks [4].

CHAPTER – 3: SYSTEM ANALYSIS

3.1 System Analysis

System analysis is a detailed examination of a system or a proposed project to identify its components, their interrelationships, and how they work together to achieve a particular goal or solve a problem. The primary purpose of system analysis is to improve the efficiency and effectiveness of the system, ensuring that it meets the needs of its users.

3.1.1 Requirement Analysis

Requirement analysis is the systematic process of gathering, documenting, and analyzing user needs and expectations for a system or project. It involves understanding the functionalities, constraints, and objectives to ensure that the resulting solution effectively meets the identified requirements.

i. Functional requirement

The functional requirements of the system include

- An actor and text-to-image generator are present.
- User provides a text prompt which leads to a series of other activities followed by visualizing the image finally.
- The system trains the model.

Actors:

User: Interacts with text to image generator.

System: Executes the model and handles processing.

Use cases:

Provide input prompt:

- User gives input text for processing.
- Includes: Validation of input, generate image, Model, Return generated image
- Extends: Display input error (for invalid input).

Visualize image:

- The user views the generated image.
- Includes: Return generated image.

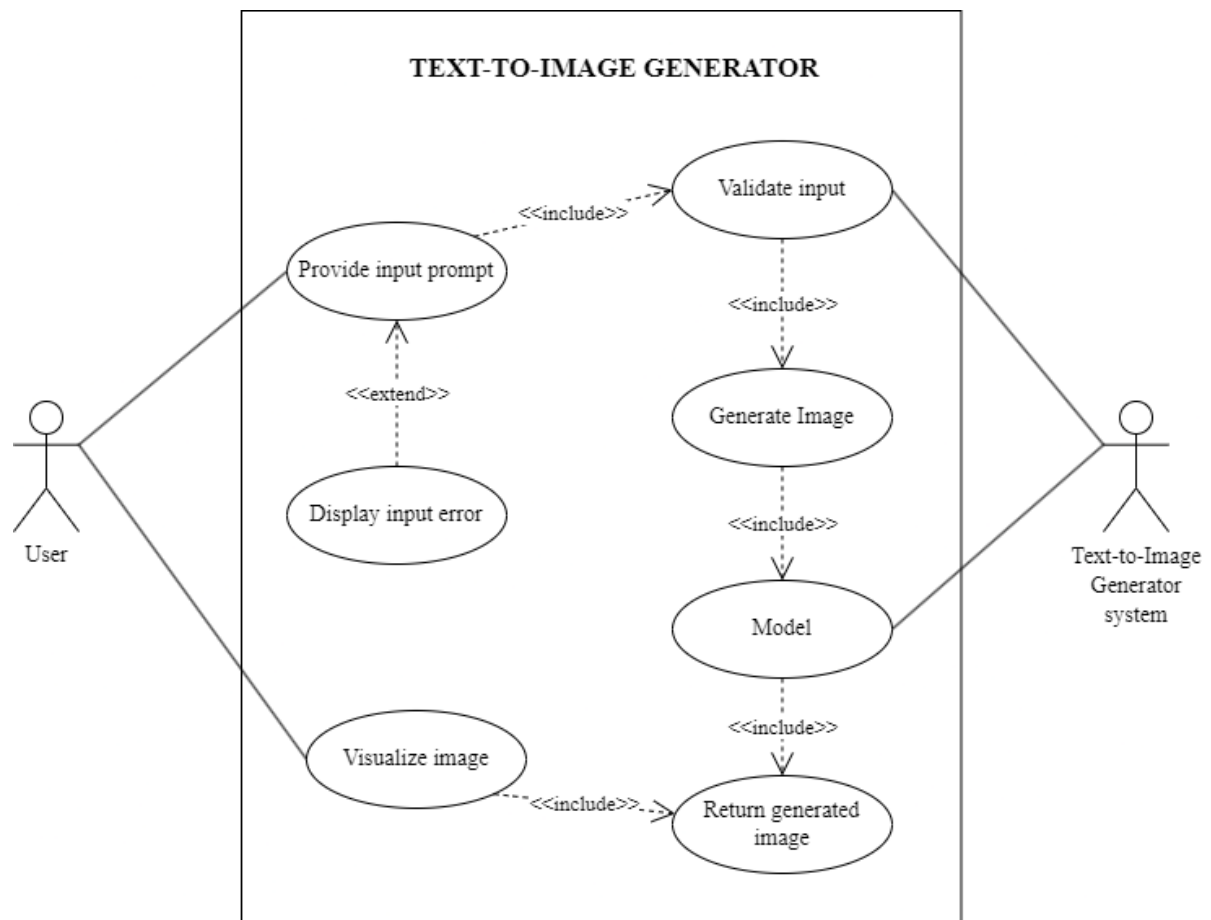


Figure 1: Use Case Diagram

ii. Non-functional requirement

The non-functional requirements of the system include:

- **Performance:** The system should be able to generate high-quality images from textual prompts within a reasonable time frame.
- **Reliability:** The model must be able to demonstrate stability by producing consistent and reliable results, minimizing the occurrence of image generation errors or inconsistencies.
- **Usability:** The user interface should be easy to use and understand. It should be intuitive and responsive, allowing users to interact with the system.

3.1.2 Feasibility Analysis

i. Technical feasibility

This project is carried out in Visual Studio Code and also in Google Colab and Kaggle Notebooks.

ii. Schedule feasibility

With the necessary resources including hardware, software, and personnel, the project can be completed within a three-month timeframe. A Gantt chart in the figure below provides a tentative schedule outlining the project's timeline.

PROCESS	MONTH 1				MONTH 2				MONTH 3			
	W1	W2	W3	W4	W1	W2	W3	W4	W1	W2	W3	W4
Research and Planning												
Model Building												
Frontend Development												
Backend development												
Model Tuning												
Documentation												

Figure 2: Gantt-Chart of the project

iii. Economic feasibility

This project is flexible in terms of economic expenses. All the software tools used for developing the project are free. The project was developed with the help of the personal devices of the team members and no other computational infrastructure was bought, rented or used throughout the project

3.1.3 Analysis

Analysis in a flow diagram refers to the examination and representation of processes, activities, or steps within a system or a project. It involves breaking down a complex system into manageable components and illustrating the sequential flow of actions or data between these components.

i. Flow diagram

The flow diagram consists of various software process stages. It is initiated by loading the dataset, followed by data cleaning and data transformation activities. Thereafter, the data is fed into the Model for training. The model's effectiveness is evaluated by validating it against the test dataset. Finally, after the training and testing process, users can generate images by providing text prompts.

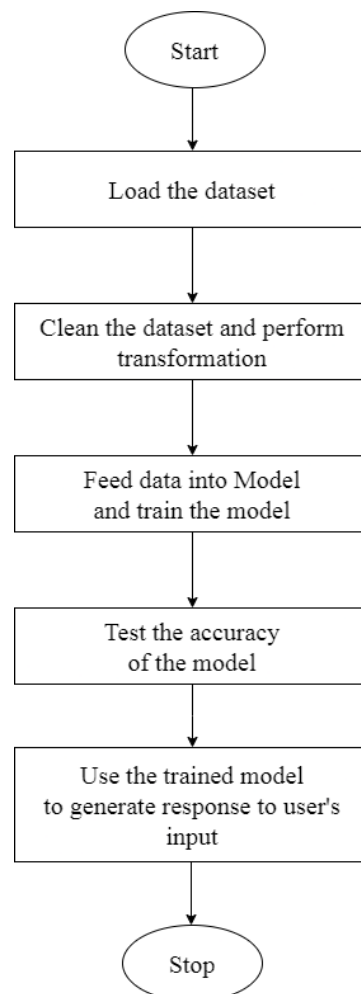


Figure 3: Flow Diagram of the system

CHAPTER – 4: SYSTEM DESIGN

4.1 Design

4.1.1 Sequence diagram

The sequence diagram in the system provides a sequential flow of activities between the user, the system and the model. In the below sequence diagram, the user initiates by entering a text prompt into the system. Then the system validates if the prompt is okay or not and delivers the appropriate message accordingly. If the prompt is valid, the system requests the model to generate the image. The model returns the required generated image to the system which then presents it to the user.

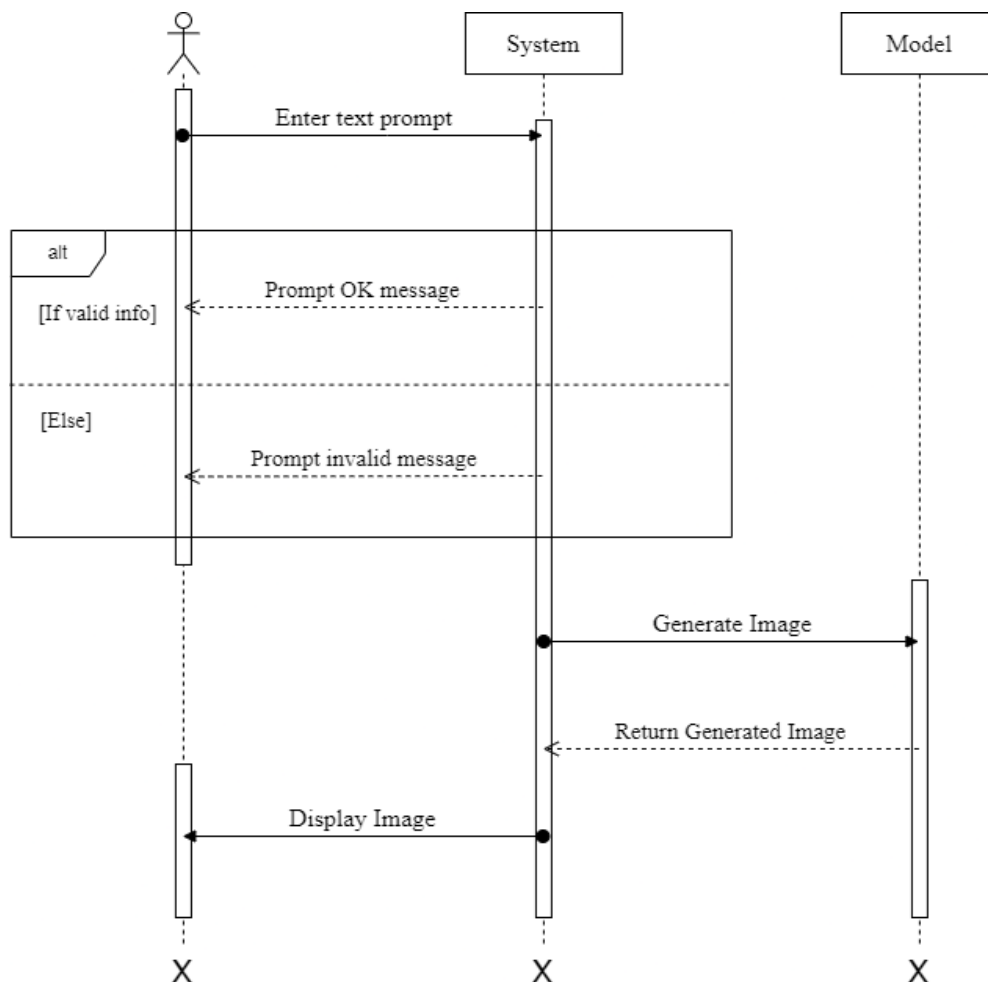


Figure 4: Sequence Diagram

4.1.2 Activity diagram

An activity diagram is the graphical workflow of all the sequential activities that take place inside the system itself. In the below activity diagram, the workflow of a system involving an actor and a text-to-image generator system is shown. The first activity is user input, where the actor enters the necessary prompt into the system. The second activity is validation, where the system checks the user input to ensure that it meets the required standards. The third activity is generating the required image using the model. The fourth activity is returning the generated image to the user for visualizing which is the last activity in the diagram.

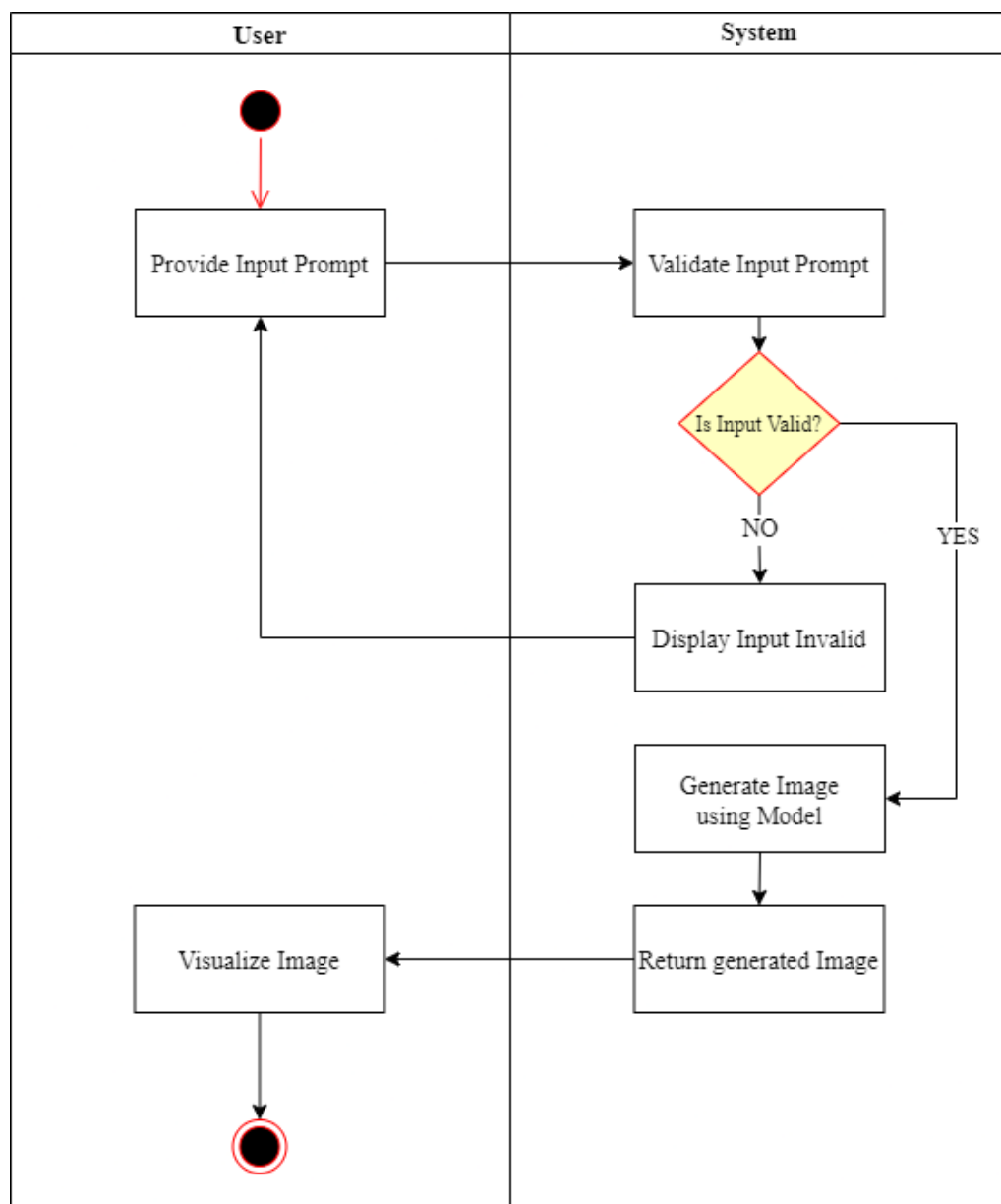


Figure 5: Activity Diagram

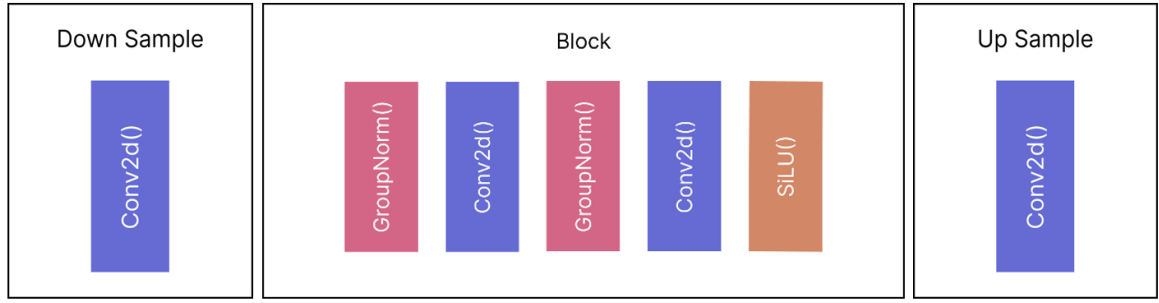


Figure 8: VAE Down Sample, VAE Block, VAE Up Sample (Left-to-right)

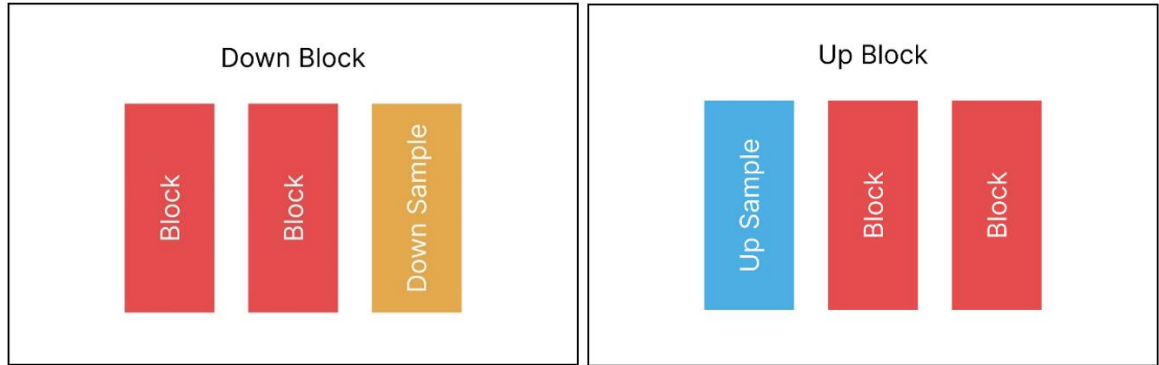


Figure 9: VAE Down and Up Block

ii. Gaussian Noise / Normal Distribution:

The Gaussian noise is added using the formula given below. The formula is simply a Normal Distribution based on mean and variance of the data.

$$q(x_t|x_{t-1}) = N(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t I)$$

iii. Conditioning:

Conditioning inputs such as text assist the model to get more knowledge about the image generation tasks which we can achieve by using text embedding. This text embedding, will be attached with time embedding and fed into the UNet.

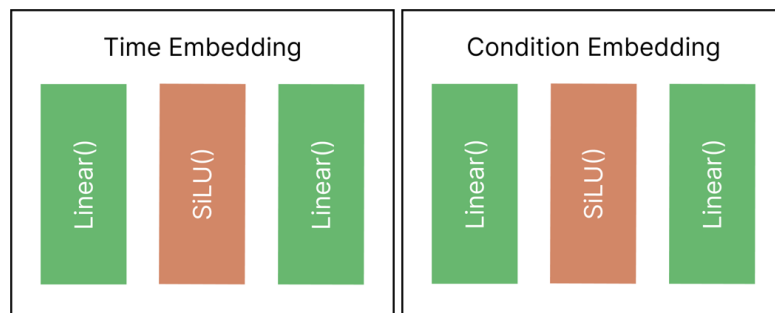


Figure 10: Time Embedding and Condition Embedding

iv. Loss function:

For the loss function, we will be using MSE loss given by the following formula:

$$L_{LDM} := E_{E(x), y, \epsilon \sim N(0,1), t} [\| \epsilon - \epsilon_{\theta}(z_t, t, \tau_{\theta}(y)) \|_2^2]$$

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (\hat{Y}_i - Y_i)^2$$

v. Denoising U-Net:

Denoising U-Net is a neural network in LDMs to predict noise of the noisy latent image. This U-Net was first introduced for image segmentation tasks but its uses are far more prevalent in other fields. This architecture will be used for predicting noise but there will be some changes in order to fuse it with LDM.

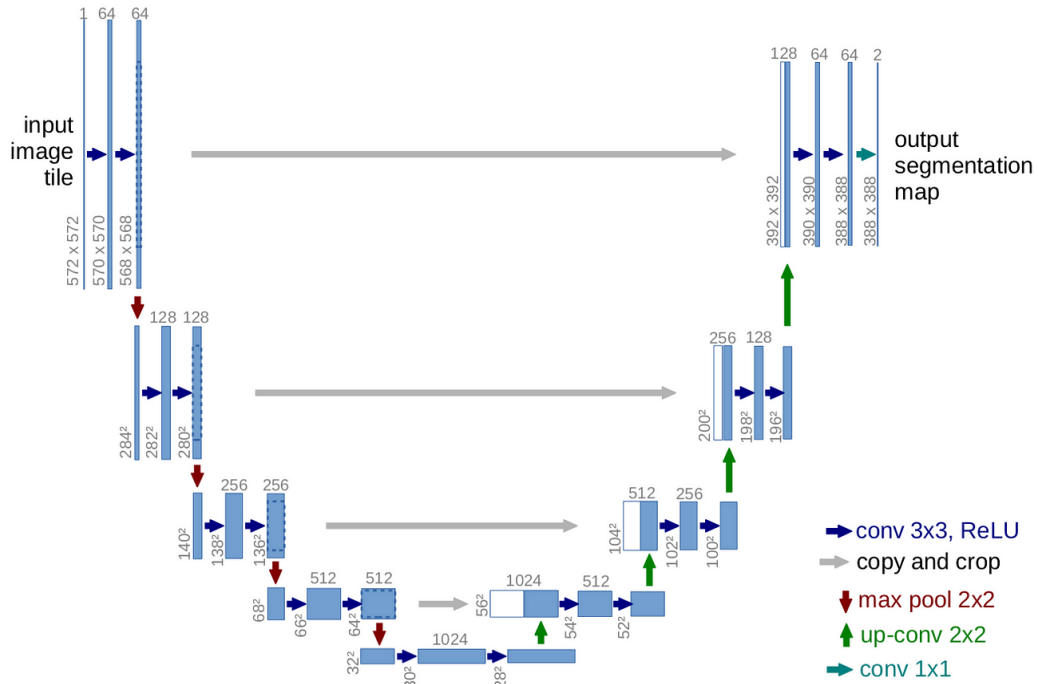


Figure 11: UNet Architecture [4]

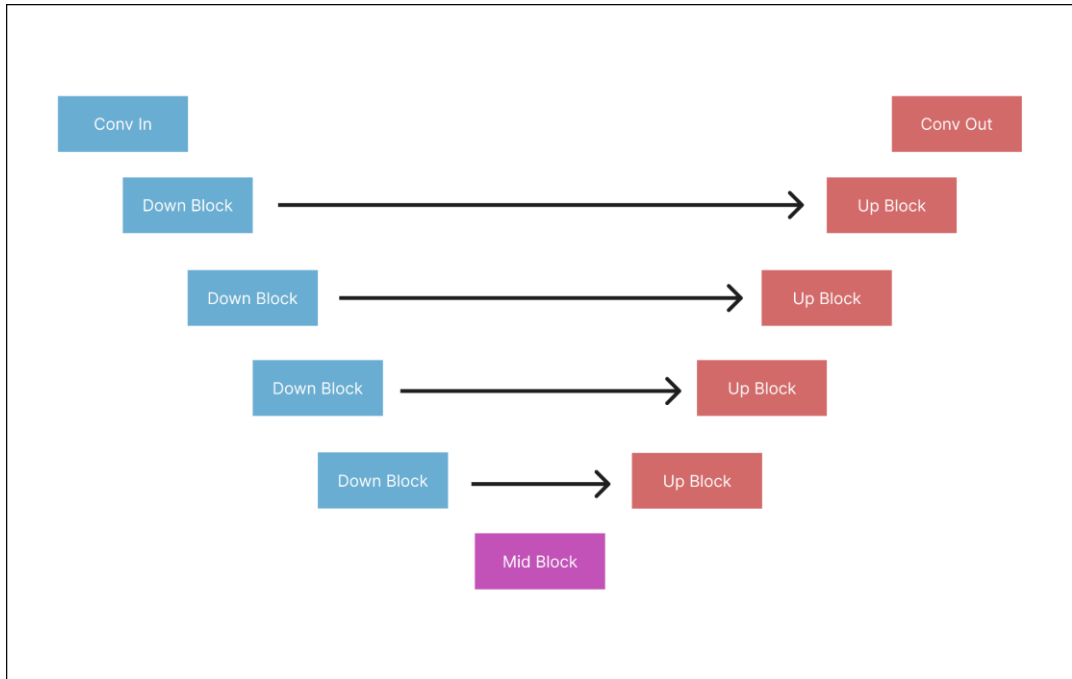


Figure 12: UNet for LDM

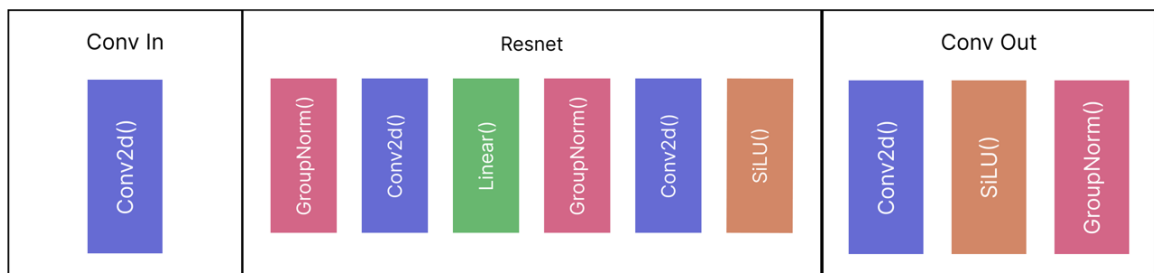


Figure 13: Conv In, Resnet, Conv Out (Left-to-right)

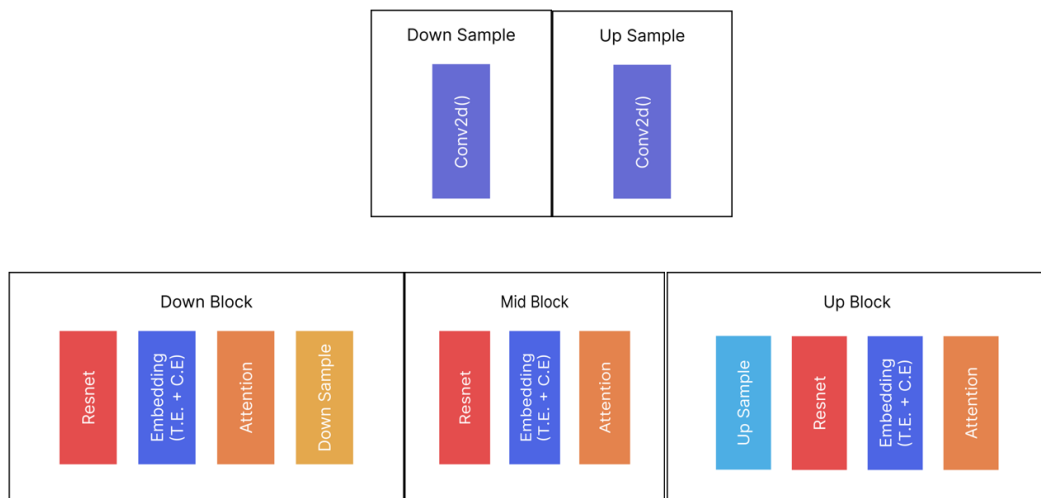


Figure 14: UNet: Down Sample, Up Sample, Down Block, Mid Block, Up Block

Training Algorithm:

Step-1: Start with an input image.

Step-2: Encode the input image into a compressed latent space using an autoencoder.

Step-3: Add Gaussian noise to the latent space to create a noisy latent space.

Step-4: Feed the noised image to a U-Net with added time embedding and condition embedding (text) to get the predicted noise.

Step-5: Remove the predicted noise from the image at timestep (T) and add noise of timestep (T-1)

Step-6: Repeat Step-4 and Step-5 until timestep (T=0).

Algorithm 1 Training	Algorithm 2 Sampling
1: repeat 2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ 3: $t \sim \text{Uniform}(\{1, \dots, T\})$ 4: $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 5: Take gradient descent step on $\nabla_{\theta} \ \epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t)\ ^2$ 6: until converged	1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 2: for $t = T, \dots, 1$ do 3: $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ if $t > 1$, else $\mathbf{z} = \mathbf{0}$ 4: $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 5: end for 6: return \mathbf{x}_0

Figure 15: Training and Sampling Algorithm [3]

Sampling Algorithm:

Step-1: Start with a random Gaussian noise.

Step-2: Feed the noised image to a U-Net with added time embedding and condition embedding (text) to get the predicted noise.

Step-3: Remove the predicted noise from the image at timestep (T) and add noise of timestep (T-1)

Step-4: Repeat Step-4 and Step-5 until a clear image is obtained.

CHAPTER – 5: IMPLEMENTATION AND TESTING

5.1 Implementation

Implementation is part of the action to put all the plans to work using different tools and programming languages.

5.1.1 Tools Used

For the implementation of this project, Python was selected as the primary programming language. In addition to Python's native functionality, several supporting libraries were utilized to facilitate the development process. The following libraries were particularly instrumental in the project's successful completion:

1. PyTorch
2. Diffusers
3. Matplotlib
4. Numpy
5. PIL library
6. Gradio

These libraries helped in achieving the desired functionality and allowed for efficient coding practices.

5.1.2 Implementation Details of Modules

i. `app.py`

This module runs the Gradio interface of the program. This launches an interface where the users interact with the program providing input text prompt and viewing the output image in Inference tab. Also, users can train their own LDM model providing different hyper parameters through the Training tab.

ii. `dataset.py`

a. CustomDataset

This class, defined in ‘dataset.py’, extends the PyTorch’s dataset class and is designed to handle the dataset of folder images. It takes the path to the image folder as input and provides methods to load and preprocess the image data.

b. NumpyDS

Similar to CustomDataset, NumpyDS is another class in 'dataset.py' that extends the CustomDataset class. It is specialized in handling latent images represented as NumPy arrays.

iii. dataloader.py:

This module contains latent_dataloader() function that simply takes the latent folder destination and loads the dataset using 'NumpyDS' class and then uses PyTorch's DataLoader to load the data.

iv. data_preprocessing.py

a. data_to_latents()

This function takes in 'VAE model', 'data_folder', 'latent_folder' and other data preprocessing parameters. Then it uses 'CustomDataset' to load the dataset and feed it into PyTorch's DataLoader which loads the dataset for preprocessing. In preprocessing, the loaded image data are encoded by the VAE and saves the latents in NumPy (.npy) file extension format in the destination folder given by 'latent_folder'.

b. data_preprocessing()

This simply calls the 'data_to_latents()' method using vae on 'data_folder' and create the latent images and put in on 'latent_folder' directory.

v. ldm.py

Here, the functionalities related to input validation and image generation are present. It initializes and manages instances of neural networks, specifically UNet and VAE, and loads their weights. A DDPM Sampler with 1000 timesteps is also instantiated for noising the image.

A 'generate_image()' function is present which involves creating random noise, converting text into tensors, passing them through UNet, and predicting noise in reverse order over a specified number of timesteps (1000). Finally after timestep reaches 0, generated image is returned.

vi. `ldm_components.py`

This module contains all the models and components used to create LDM such as VAE, DDPM and U-Net.

vii. `inference.py`

This code defines an ``inference`` function that generates an image using a DDPM and a U-Net model. The function takes several arguments, including the U-Net model, a DDPM scheduler, a VAE model, and the device. It creates a tensor of random noise, denoises the image using the U-Net model and DDPM scheduler, and saves intermediate and final images if specified. The function returns the final denoised image and a plot of saved images.

viii. `trainer.py`

This code defines two functions, ``train`` and ``cancel``. The ``train`` function trains a U-Net model using a DDPM and a VAE model. The function takes several arguments, including the number of epochs, batch size, learning rate, loss function, optimizer, device, and optional arguments for the data folder, preprocessing, and progress bar. The function initializes the DDPM scheduler, VAE model, U-Net model, and optimizer, and loads the data from a specified folder. The function then trains the U-Net model using the DDPM scheduler and optimizer, and saves the best model and logs the training loss. The ``cancel`` function sets a global flag to stop the training process.

ix. `utils.py`

This code defines several functions for data processing and visualization. The ``get_class_tensor`` function maps a class name to a tensor index, and the ``get_class_name`` function maps a tensor index to a class name. The ``latents_to_pil`` function converts a tensor of latent variables to a PIL image, and the ``plot_images_from_folder`` function plots a grid of images from a specified folder. The ``save_images`` function saves a grid of images to a specified file, and the ``plot_results`` function plots a line chart of a specified result over time. These functions are used for data processing, visualization, and evaluation in the larger context of a machine learning pipeline.

5.2.3 User Testing:

- Description: Evaluate the usability and user-friendliness of the Text-to-Image generator interface.
- Test Steps:
 - Ask users to perform specific tasks, such as entering a textual prompt, customizing image parameters, and viewing generated images.
 - Observe users' interactions with the interface, noting any difficulties, confusion, or inefficiencies.
- Expected Outcome: Users can easily navigate the interface, perform tasks intuitively, and achieve desired outcomes efficiently and easily.

5.3 Result Analysis

5.3.1 Loss Analysis:

The model converged to a loss of 0.043 i.e. 4.3% after 1000 epochs. Hence, it demonstrated a strong capability to learn from the training data and perform generation. This performance is encouraging for real-world applications, suggesting the model can effectively generate text relevant images. Also we used MSE loss as our loss function.

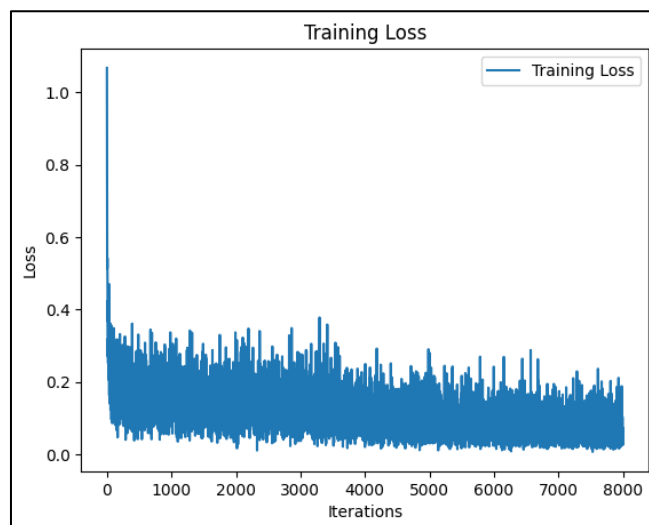


Figure 20: Loss Curve (MSE Loss)

CHAPTER – 6: CONCLUSION AND FUTURE RECOMMENDATIONS

6.1 Conclusion

In conclusion, the Text-to-Image Generator project has demonstrated significant advancements in the field of generative modeling, showcasing the potential for synthesizing visually compelling images from textual descriptions. Through the integration of state-of-the-art techniques such as the diffusion model, UNet, and VAE, along with attention mechanisms and embeddings, the project has successfully bridged the gap between natural language processing and computer vision. The utilization of PyTorch and the Diffusers library has provided a robust framework for implementation, ensuring scalability, efficiency, and real-time processing capabilities. By enabling users to creatively express their ideas through textual prompts and generating corresponding images, this project has laid the foundation for diverse applications in art, design, education, marketing, and beyond.

6.2 Future Recommendations

1. **Enhanced Attention Mechanisms:** Explore advanced attention mechanisms, such as multi-head attention or self-attention mechanisms, to further improve the model's understanding of textual input and enhance the relevance of generated images.
2. **Integration of Transfer Learning:** Investigate the integration of transfer learning techniques to leverage pre-trained models and further enhance the quality and diversity of generated images, particularly in scenarios with limited training data.
3. **User Interface Improvements:** Develop a user-friendly interface that allows users to interact with the system intuitively, providing options for customization, exploration of generated images, and feedback mechanisms for continuous improvement.
4. **Ethical Considerations:** Conduct thorough analyses to identify and mitigate potential biases in the training data and generated images, ensuring that the system adheres to ethical standards and promotes diversity and inclusivity.
5. **Scalability and Optimization:** Optimize the system for scalability and efficiency, exploring techniques such as distributed training, model compression, and hardware acceleration to handle larger datasets and increase processing speed.

REFERENCES

- [1] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, Bjorn Ommer, “High-Resolution Image Synthesis with Latent Diffusion Models”, *Arxiv*, Apr. 13, 2022. [Online]. Available: <https://arxiv.org/pdf/2112.10752.pdf>
- [2] Patrick Esser, Robin Rombach, Bjorn Ommer, “Taming Transformers for High-Resolution Image Synthesis” *Arxiv*, Dec. 16, 2020. [Online]. Available: <https://arxiv.org/pdf/2006.11239.pdf>
- [3] Jonathan Ho, Ajay Jain and Pieter Abbeel, “Denoising Diffusion Probabilistic Models”, *Arxiv*, Dec. 16, 2020. [Online]. Available: <https://arxiv.org/pdf/2006.11239.pdf>
- [4] Olaf Ronneberger, Philipp Fischer and Thomas Brox, “U-Net: Convolutional Networks for Biomedical Image Segmentation”, *Arxiv*, May. 18, 2015. [Online]. Available: <https://arxiv.org/pdf/1505.04597.pdf>
- [5] Alex Nichol and Prafulla Dhariwal, “Improved Denoising Diffusion Probabilistic Models,” Feb. 2021. [Online]. Available: <https://arxiv.org/pdf/2102.09672.pdf>
- [6] J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, S. Ganguli, and S. Edu, “Deep Unsupervised Learning using Nonequilibrium Thermodynamics,” Nov. 2015. Available: <https://arxiv.org/pdf/1503.03585.pdf>
- [7] D. P. Kingma and M. Welling, “Auto-Encoding Variational Bayes,” arXiv.org, Dec. 20, 2013. Available: <https://arxiv.org/abs/1312.6114>

APPENDIX-I

app.py:

```
import gradio as gr
from ldm import generate_image
from trainer import train, cancel

with gr.Blocks(title="Text-to-Image Generator") as app:
    gr.Markdown("# Text-to-Image Generator: 7th Semester Project")
    with gr.Tab("Inference"):
        gr.Markdown("### Insert a text and generate a face.")
        with gr.Row():
            with gr.Column():
                text = gr.Textbox(label="Enter text prompt")
                gr.Examples([
                    ["Barack Obama"],
                    ["Cristiano Ronaldo"],
                    ["Donald Trump"]
                ], inputs=[text])
                generate_btn = gr.Button("Generate")
                gr.ClearButton([text])

            with gr.Column():
                image_output = gr.Image(label="Generated image", width=500)

            with gr.Column():
                gr.Image(label="Model Loss", value="./images/final.png", width=500)

        with gr.Row():
            with gr.Column():
                inference_output = gr.Image(label="Timesteps")
                generate_btn.click(
                    generate_image,
                    inputs=[text],
                    outputs=[image_output, inference_output],
                )

    with gr.Tab("Training"):
        gr.Markdown("### Train a Text-to-Image Generator by inserting following Hyperparameters.")
        with gr.Row():
            epochs = gr.Textbox(label="Epochs", value=1000)
            batch_size = gr.Textbox(label="Batch Size", value=20)
```

```

lr = gr.Textbox(label="Learning Rate", value=0.001)

with gr.Row():
    loss_fn = gr.Radio(["MSE", "MAE"], label="Loss Function", value="MSE")
    optimizer = gr.Radio(["Adam", "SGD"], label="Optimizer", value="Adam")
    latent_folder = gr.Textbox(label="Latent Folder Location",
value="./data/face/images")

with gr.Row():
    preprocess = gr.Radio(["Yes", "No"], label="Preprocessing", value="No")
    data_folder = gr.Textbox(label="Data Folder Location")

with gr.Row():
    train_btn = gr.Button("Train")
    cancel_btn = gr.Button("Cancel")

progress_output = gr.Textbox()

train_btn.click(
    train,
    inputs=[epochs, batch_size, lr, loss_fn, optimizer, latent_folder, preprocess,
data_folder],
    outputs=progress_output
)
cancel_btn.click(cancel)

with gr.Tab("Scratch Results"):
    gr.Markdown("### Results Generated by Scratch Model.")
    with gr.Column():
        gr.Image(value="./images/scratch-1.png")
        gr.Image(value="./images/scratch-2.png")

if __name__ == "__main__":
    app.launch()

```

dataset.py:

```

import os
from PIL import Image
import numpy as np
from torch.utils.data import Dataset

class CustomDataset(Dataset):
    def __init__(self, root_dir, transform=None):

```



```

self.root_dir = root_dir
self.transform = transform
self.classes = os.listdir(root_dir)
self.class_to_idx = {cls: i for i, cls in enumerate(self.classes)}
self.file_list = []

for cls in self.classes:
    class_path = os.path.join(self.root_dir, cls)
    files = os.listdir(class_path)
    self.file_list.extend([(cls, os.path.join(class_path, file)) for file in files])

def __len__(self):
    return len(self.file_list)

def __getitem__(self, idx):
    class_name, img_path = self.file_list[idx]
    image = Image.open(img_path).convert("RGB").resize((256, 256))

    label = self.class_to_idx[class_name]

    if self.transform:
        image = self.transform(image)

    return image, label

class NumpyDS(CustomDataset):
    def __getitem__(self, idx):
        class_name, img_path = self.file_list[idx]

        image = np.load(img_path)
        label = self.class_to_idx[class_name]

        return image, label

```

dataloader.py:

```

from dataset import NumpyDS
from torch.utils.data import DataLoader

def latent_dataloader(latent_folder, batch_size, shuffle):
    data_set = NumpyDS(root_dir=latent_folder)
    return DataLoader(dataset=data_set, batch_size=batch_size, shuffle=shuffle)

```

data_preprocessing.py:

```
import os
import numpy as np
from torchvision import transforms
from torch.utils.data import DataLoader

from dataset import CustomDataset

def data_to_latents(model, data_folder, latent_folder, transform, batch_size):
    data_set = CustomDataset(root_dir=data_folder, transform=transform)
    label = data_set.classes
    data_loader = DataLoader(dataset=data_set, batch_size=batch_size)
    name = 0
    for X, y in data_loader:
        X_encode = model.encode(X.to(model.device)).latent_dist.sample() * 0.18215
        for X_encode_i, y_i in zip(X_encode, y):
            y_i_label = latent_folder / label[y_i.item()]
            y_i_label.mkdir(parents=True, exist_ok=True)
            file_path = os.path.join(y_i_label, f"{name}.numpy")
            name += 1
            np.save(file_path, X_encode_i.cpu().numpy())

def data_preprocessing(vae, data_folder, latent_folder):
    transform = transforms.Compose([
        transforms.ToTensor(),
        lambda x: x * 2.0 - 1.0
    ])
    data_to_latents(vae, data_folder, latent_folder, transform=transform, batch_size=5)
```

ldm.py:

```
import torch

import ldm_components
from inference import inference

device = torch.device("cuda" if torch.cuda.is_available() else "cpu")

num_of_generation = 1

ddpm_scheduler = ldm_components.ddpm_scheduler()
vae = ldm_components.vae().to(device)
UNET = ldm_components.unet().to(device)

UNET.load_state_dict(torch.load("./saved/face_unet.pth", map_location=device))
```

```
def generate_image(text):
    if text == "": return "./images/error.jpg"
    return inference(unet, ddpm_scheduler, vae, device, text=text)
```

ldm_components.py:

```
from diffusers import AutoencoderKL, DDPM Scheduler, UNet2DModel
```

```
def vae():
    return AutoencoderKL.from_pretrained("stabilityai/sd-vae-ft-ema").requires_grad_(False)
```

```
def ddpm_scheduler():
    return DDPM Scheduler(num_train_timesteps=1000, beta_start=0.0001,
beta_end=0.02)
```

```
def unet():
    return UNet2DModel(
        sample_size=32,
        in_channels=4,
        out_channels=4,
        layers_per_block=2,
        block_out_channels=(32, 64, 128, 256),
        down_block_types=(
            "DownBlock2D",
            "AttnDownBlock2D",
            "AttnDownBlock2D",
            "AttnDownBlock2D"
        ),
        up_block_types=(
            "AttnUpBlock2D",
            "AttnUpBlock2D",
            "AttnUpBlock2D",
            "UpBlock2D"
        ),
        class_embed_type="timestep"
    )
```

inference.py:

```
import torch
from utils import get_class_tensor, latents_to_pil, save_images, plot_images_from_folder
```

```
def inference(unet, ddpm_scheduler, vae, device, num_images=1, text=None,
img_name=None):
```

```

noisy_images = torch.randn((num_images, 4, 32, 32)).to(device)
if text:
    label = get_class_tensor(text)
    if label is None:
        return "./images/error.jpg"
    labels = label.expand(num_images).to(device)
else:
    labels = torch.randint(3, size=(noisy_images[:num_images].shape[0],)).to(device)

for timestep in ddpm_scheduler.timesteps:
    unet.eval()
    img_checkpoints = [999, 700, 500, 200, 150, 120, 100, 80, 60, 40, 20, 0]
    with torch.inference_mode():
        noise = unet(noisy_images, timestep, labels).sample
        noisy_images = ddpm_scheduler.step(noise, timestep, noisy_images).prev_sample
        if timestep.item() in img_checkpoints and text:
            latents_to_pil(vae, noisy_images,
save_path=f"./images/inference/{timestep}.png")

            images = latents_to_pil(vae, noisy_images)

if text:
    inference_img = plot_images_from_folder("./images/inference")
    return images[0], inference_img

save_images(img_name, images, labels)

```

trainer.py:

```

import torch
from torch import nn

import gradio as gr

from pathlib import Path

import ldm_components
from inference import inference
from data_processing import data_preprocessing
from utils import plot_results
from dataloader import latent_dataloader

continue_training = True

def train(

```

```

epochs,
batch_size,
lr,
loss_fn,
optimizer,
device,
latent_folder=None,
preprocess=False,
data_folder=None,
progress=gr.Progress()
):
    device = torch.device("cuda" if torch.cuda.is_available() else "cpu")

    ddpm_scheduler = ldm_components.ddpm_scheduler()
    vae = ldm_components.vae().to((device))
    unet = ldm_components.unet().to(device)

    loss_fn = nn.MSELoss() if loss_fn == "MSELoss" else nn.L1Loss()
    if optimizer == "Adam":
        optimizer = torch.optim.AdamW(unet.parameters(), lr=float(lr))
    else:
        optimizer = torch.optim.SGD(unet.parameters(), lr=float(lr), momentum=0.999,
weight_decay=0.0001)

    dataloader = latent_dataloader(Path("./data/face/latents/"), batch_size=int(batch_size),
shuffle=True)

    if preprocess and latent_folder: data_preprocessing(vae, data_folder, latent_folder)

    losses = []

    best_loss = float("inf")
    best_model_path = "./saved/best_model.pth"

    print_loss = 0
    global continue_training
    for epoch in range(int(epochs)):
        train_loss = 0

        unet.train()
        for latent_images, labels in progress.tqdm(dataloader, desc=f"Training: Epoch:
{epoch+1} | Training Loss: {print_loss:.3f} |"):
            if not continue_training: break
            latent_images, labels = latent_images.to(device), labels.to(device)

```

```

        timesteps = torch.randint(0, ddpm_scheduler.config.num_train_timesteps,
(latent_images.shape[0], )).to(device)
        noise = torch.randn(latent_images.shape).to(device)
        noisy_images = ddpm_scheduler.add_noise(latent_images, noise, timesteps)

        pred_noise = unet(noisy_images, timesteps, labels).sample
        loss = loss_fn(pred_noise, noise)
        train_loss += loss.item()
        losses.append(train_loss)
        optimizer.zero_grad()
        loss.backward()
        optimizer.step()

    if not continue_training: break

    if train_loss < best_loss:
        best_loss = train_loss
        torch.save(unet.state_dict(), best_model_path)

    if epoch % 10 == 0:
        inference(unet, ddpm_scheduler, vae, device, num_images=5, img_name=epoch)

    print_loss = train_loss / len(dataloader)
    print(f"Epoch: {epoch} | Loss: {print_loss:.3f}")

    plot_results(losses, "Training Loss")
    return "Training Complete" if continue_training else "Training Cancelled"

def cancel():
    global continue_training
    continue_training = False

```

utils.py:

```

import os
import re

import torch
import matplotlib.pyplot as plt
from PIL import Image

def get_class_tensor(class_name):
    class_mapping = {"barack obama": 0, "cristiano ronaldo": 1, "donald trump": 2}

```

```

lower_class_name = class_name.lower()
if lower_class_name in class_mapping:
    class_idx = class_mapping[lower_class_name]
    return torch.tensor(class_idx)
return None

def get_class_name(class_idx):
    class_mapping = {0: "barack obama", 1: "cristiano ronaldo", 2: "donald trump"}
    if class_idx in class_mapping:
        return class_mapping[class_idx]
    return None

def latents_to_pil(vae, latents, save_path=None):
    latents = (1 / 0.18215) * latents
    with torch.inference_mode():
        image = vae.decode(latents).sample
    image = (image / 2 + 0.5).clamp(0, 1)
    image = image.detach().cpu().permute(0, 2, 3, 1).numpy()
    images = (image * 255).round().astype("uint8")
    pil_images = [Image.fromarray(image) for image in images]
    if save_path:
        pil_images[0].save(save_path)
    return pil_images

def plot_images_from_folder(folder_path):
    save_path = "./images/inference_img.png"
    _, axes = plt.subplots(2, 6, figsize=(20, 8))
    image_files = sorted(
        os.listdir(folder_path),
        key=lambda x: int(re.findall(r"\d+", x)[0]),
        reverse=True
    )
    for i, image_file in enumerate(image_files):
        image_path = os.path.join(folder_path, image_file)
        image = Image.open(image_path)
        ax = axes[i // 6, i % 6]
        ax.imshow(image)
        ax.set_title(f"t={image_file}")
        ax.axis('off')
    plt.tight_layout()
    plt.savefig(save_path)
    plt.close()
    return Image.open(save_path)

```

```

def save_images(img_name, images, labels):
    if not os.path.exists("results"):
        os.makedirs("results")
    fig = plt.figure(figsize=(15, 3))
    for i in range(5):
        ax = fig.add_subplot(1, 5, i+1)
        ax.set_title(get_class_name(labels[i].item()))
        ax.imshow(images[i])
        ax.axis("off")
        plt.savefig(f'./images/output.png')
    plt.close()

```

```

def plot_results(result, title):
    batches = range(1, len(result) + 1)
    plt.plot(batches, result, label=title)
    plt.title(title)
    plt.xlabel("Iterations")
    plt.ylabel("Loss")
    plt.legend()
    # plt.show()
    plt.savefig(f'./images/error_train.png')
    plt.close()

```