# WHIMCARNATOR - SKETCH TO REAL IMAGE TRANSFORMER



By

**MUZAMMIL GHAFFAR (LEADER)**
**HUSNAIN MUSTAFA**
**MUHAMMAD USAMA BUTT**
**SUBHAN UD DIN BAJWA**

Supervised by:

**DR. NAUMAN ALI KHAN**

Submitted to the faculty of Department of Computer Software Engineering,

Military College of Signals, National University of Sciences and Technology,

Islamabad,

in partial fulfilment for the requirements of B.E Degree in Software Engineering.

June 2023

In the name of ALLAH, the Most benevolent, the Most Courteous

# DECLARATION OF ORIGINALITY

We hereby declare that no portion of work presented in this thesis has been submitted in support of another award or qualification in either this institute or anywhere else.

# ACKNOWLEDGEMENTS

# Table of Contents

# LIST OF FIGURES

# ABSTRACT

Whimcarnator is a sketch-to-image transformer application that aims to generate photorealistic images from user-drawn sketches or uploaded images. This project addresses the increasing popularity of touchscreens and the need for sketch-based applications in various fields such as graphics and game design, e-commerce, law enforcement, and accessibility tools. The purpose of Whimcarnator is to provide a platform that lets users generate photorealistic images from their sketches, which can be used in various ways, such as identifying suspects, searching for products online, capturing scenes or objects, and recognizing patterns in vague images. Additionally, Whimcarnator allows other applications to integrate with it, providing them with the capability to retrieve images and create photorealistic images. This document presents the software requirements specification for the Whimcarnator project, outlining the features and requirements of the application for developers and prospective clients alike.

# 1. INTRODUCTION

## 1.1. Overview

Sketches have been an essential means of communication since prehistoric times. However, in recent years, the advent of touchscreens and digital sketching devices has made the creation of sketches much easier, and consequently, sketch-oriented applications have become increasingly popular. Sketches are a powerful and natural visual representation of the world around us. They are highly illustrative and can be used to depict objects, stories, or concepts. With the abundance of touchscreens and touchpads, we can now incorporate sketches in our digital systems to tap into their potential.

One such system is Whimcarnator, a sketch-to-image transformer that generates photorealistic images from sketches. Whimcarnator is built using a pre-trained model and a conditional fine-tuning method. The pre-trained model is trained on a large dataset of images and has learned to extract meaningful features from images. The conditional fine-tuning method allows us to fine-tune the pre-trained model to generate images conditioned on the input sketch.

The purpose of Whimcarnator is to take a sketch from a user or allow the user to create a sketch on the provided canvas and generate variations of respective photorealistic images. These images can be used in several ways, including product search on e-commerce sites, quickly capturing a scene or object after drawing a sketch of it, creating models in graphics and game designing, as an accessibility tool to help users recognize patterns and vague images, or helping law-enforcement agencies to identify suspects by generating their photorealistic images from sketches.

In order to develop the Whimcarnator system, we will employ a generative adversarial network (GAN) trained on an extensive collection of images. The GAN will be composed of a generator and a discriminator. The generator's task is to create photorealistic images based on the input sketch, while the discriminator assesses the realism of the generated images. To achieve this, the generator will undergo fine-tuning using a conditional approach, wherein the input sketch serves as a guide for the generator to produce photorealistic images that accurately represent the input sketch.

In addition, we will employ methods like adversarial diffusion up-sampling and classifier-free normalized guidance to improve the synthesis quality of the produced images. These approaches are designed to address challenges like accurately aligning with provided inputs and preventing small objects from being omitted in the generated images.

Overall, Whimcarnator is an innovative application that leverages the power of touchscreens and digital sketching devices to generate photorealistic images from sketches. It has several potential applications and can be used to facilitate various tasks, from e-commerce product search to law enforcement investigations.

## 1.2. Scope

The scope of the project includes the development of a user-friendly web or mobile application that provides a canvas for the user to draw sketches and then generates photorealistic images from those sketches. The application will also allow the user to upload an existing sketch and generate a photorealistic image from it.

Furthermore, the project aims to extend the functionality of the proposed model to integrate it with other applications, such as e-commerce platforms, graphic designing tools, and law enforcement agencies. The integration with e-commerce platforms will enable users to search for products using sketches, while the integration with graphic designing tools will allow designers to create photorealistic images of their models. The integration with law enforcement agencies will help them generate photorealistic images of suspects from sketches, aiding them in investigations.

Overall, the scope of the project is to develop a comprehensive solution that addresses the need for generating photorealistic images from sketches, and provides a platform for its integration with other applications, thereby enhancing its utility and impact.

## 1.3. Working Principle

Whimcarnator is a web application built using Next.js, a popular framework for building React applications. The application allows users to create freehand sketches on a provided canvas, and using a machine learning backend built in Python, generates photorealistic images based on the sketch.

The working principle of Whimcarnator is to generate a photorealistic image from a user-provided sketch using image-to-image translation techniques. The process involves training a deep neural network on a large dataset of paired sketches and corresponding real images. The trained network can then take a user's sketch as input and output a photorealistic image that closely matches the input sketch.

## 1.4. Model Training

To train the deep learning model for Whimcarnator, we used a modified version of the pre-training method introduced in the paper "Pretrained Image Transformers" by Wang et al. The architecture used is a variant of the U-Net model, which has been widely used in image-to-image translation tasks. The model is trained on a large dataset of paired sketches and corresponding photorealistic images. We used the Sketchy dataset, which consists of more than 125,000 sketches across 125 categories, and the COCO-Stuff dataset, which contains over 10,000 images with dense object-level annotations, as our training data. The model is trained using the adversarial training method, which includes a discriminator network that distinguishes between the generated photorealistic images and the ground truth images.

To train the model, we used a high-performance GPU cluster to speed up the training process. We used the PyTorch deep learning framework to implement the model and the training process. The training process involved several stages, starting with pre-training the U-Net model on the Sketchy dataset, followed by fine-tuning on the COCO-Stuff dataset using the adversarial training method. We used various techniques such as data augmentation, dropout regularization, and learning rate scheduling to prevent overfitting and improve the generalization of the model.

The training process took several days to complete, but the resulting model is capable of generating high-quality photorealistic images from sketches in real-time. The trained model will be deployed on the machine learning backend, which will be responsible for generating photorealistic images from the user's sketches in real-time.

## 1.5.   Relevant Sustainable Development Goals

There are several relevant sustainable development goals (SDGs) that can be associated with the Whimcarnator project:

**SDG 9 - Industry, Innovation, and Infrastructure:** This project aims to develop an innovative tool that utilizes touchscreens and sketching to create photorealistic images. By developing such a tool, the project contributes to the development of sustainable infrastructure and encourages innovation.

**SDG 12 - Responsible Consumption and Production:** The Whimcarnator project aims to enable users to create photorealistic images from their sketches. By doing so, the project aims to reduce the reliance on traditional art supplies, which may have negative impacts on the environment.

**SDG 13 - Climate Action:** The project aims to reduce the carbon footprint associated with the creation of images by using digital means rather than traditional art supplies. By doing so, the project supports efforts to reduce greenhouse gas emissions and mitigate the impacts of climate change.

**SDG 4 - Quality Education:** The project has the potential to support quality education by enabling students to create photorealistic images in an efficient and sustainable manner. By doing so, the project may help to enhance the learning experience and promote creativity in the classroom.

**SDG 8: Decent Work and Economic Growth:** The Whimcarnator project has the potential to support economic growth by providing a tool that can be used in a variety of industries, including art, graphic design, and e-commerce. By doing so, the project may help to create new job opportunities and support the growth of these industries.

## 1.6. Organization of report

**Chapter 2: Literature Review**

- Introduction

- Existing Solutions

- Related Work

- Drawbacks of Existing Solutions

- Conclusion

**Chapter 3: System Features**

- Deliverables

- Extended Scope

**Chapter 4: Non-functional Requirements**

- Performance Requirements

- Safety Requirements

- Security Requirements

- Software Quality Attributes

**Chapter 5: Design and Development**

- Introduction

- Design

- Architecture

**Chapter 6: Architectural Design**

- Architecture Design

- Decomposition Diagram

**Chapter 7: System Requirements Specifications**

- External Interface Requirements

- Software Interfaces

- Communication Interfaces

**Chapter 8: Implementation and Testing**

- Evaluation metrics and methodology

- Performance analysis of the model

- User testing and feedback

- Analysis of challenges faced during development

**Chapter 9: Conclusion**

- Summary of the project

- Achievements and limitations

- Conclusions drawn from the results

- Implications for future work

**Chapter 10: Future Work**

- Further improvements to the project

- Potential commercialization opportunities

- Suggestions for future research and development

**Chapter 11: References**

- List of sources referenced used to develop the project.

# 2. LITERATURE OVERVIEW

## 2.1. Introduction

This chapter presents a literature review of relevant research studies and articles related to the project Whimcarnator. The chapter aims to identify the gaps in the existing literature and provide a comprehensive understanding of the state-of-the-art technologies used in the project.

## 2.2. Existing solutions

In recent years, various approaches have been developed for image-to-image translation tasks. However, these approaches have several drawbacks that limit their effectiveness in addressing the challenges posed by this problem. In this chapter, we will review some of the existing solutions and their drawbacks.

### 2.2.1. Conditional Generative Adversarial Networks (cGANs)

Among the most prevalent techniques for image-to-image conversion is the use of conditional generative adversarial networks (cGANs). These networks depend on a discriminator to identify discrepancies with authentic images. cGANs have proven effective in numerous image-to-image translation assignments, including image colorization, super-resolution, and style transfer. Nonetheless, cGANs are susceptible to instability and mode collapse, potentially resulting in subpar output quality.

### 2.2.2. Autoregressive Models

Autoregressive models, which leverage the exceptional expressivity of transformers, have demonstrated encouraging outcomes in image-to-image translation endeavors. Despite this, they suffer from slow inference speeds and a tendency to overfit. Moreover, they handle different tasks individually and rely on restricted task-specific training data, potentially hindering their overall efficacy.

### 2.2.3. Multi-Task Learning

Various research initiatives strive to develop a comprehensive model for a range of translation tasks through the use of multi-task learning.. However, these methods still require significant amounts of task-specific training data and may not be effective in handling complex and diverse translation tasks.

### 2.2.4. Pretraining

Pretraining is essential for modern vision tasks, but the majority of pretraining methods concentrate on discriminative tasks, while pretraining for visual synthesis remains underexplored. Previous efforts have attempted to use pretrained models as generative priors for tasks such as conditional image synthesis, image editing, and restoration. These approaches frequently employ a GAN latent space, which encodes input semantics and facilitates significant manipulation. However, GANs excel primarily in representing specific image categories and struggle with mode dropping and stability concerns, rendering them inadequate as generative priors for general images.

### 2.2.5. Diffusion Models

Diffusion and score-based models have recently emerged to show competitive generation quality across various benchmarks. Diffusion models have demonstrated extraordinary capacity when trained with large-scale text-image pairs. However, Saharia et al. demonstrate the potential of using the diffusion model for image-to-image translation, but they only show results for data-rich problems, e.g., image colorization. Diffusion models are insufficient to serve as generative priors for general images.

## 2.3.  Related Work

Several research studies have explored the task of sketch-to-image translation using GANs. One notable work is SketchyGAN, which proposed a GAN-based model for generating realistic images from sketches. The model consisted of a generator network and a discriminator network, both of which were trained on a large dataset of sketches and corresponding real images. Another study proposed a conditional GAN-based model for sketch-to-image translation, where the generator was conditioned on a specific object class to generate images of that class from the input sketch.

## 2.4.  Drawbacks of Existing Solutions

Despite significant research efforts, the existing solutions for image-to-image translation still have several drawbacks, including:
- Prone to instability and mode collapse
- Slow to inference and prone to overfitting
- Require significant amounts of task-specific training data

- Inadequate for functioning as generative priors for a broad range of images

In order to tackle these limitations, we introduce a new method known as Whimcarnator, which utilizes a thoroughly pretrained diffusion model as a universal generative prior, enabling it to handle multiple synthesis tasks without the need for task-specific adjustments or hyperparameter optimization.

## 2.5. Conclusion

This chapter provided a literature review of relevant research studies and articles related to the project Whimcarnator. The chapter discussed the state-of-the-art technologies used in the project, including sketch-to-image translation, GANs, and pretrained models. The review highlighted the gaps in the existing literature and identified the need for further research in this area.

# 3. SYSTEM FEATURES

## 3.1.  Deliverables

Below are the requirements that will be expected to be delivered:

### 3.1.1.  Canvas

The system will provide an interface where users will be able to draw their sketches. This interface will include the following characteristics.

- Users can draw freely on all portions of the canvas.
- Users can select different brush sizes varying in width.

### 3.1.2.  Image Generation

The system should generate a variation of photorealistic images from the provided sketch.

## 3.2.  Extended Scope

The extended scope enlists all the FRs that if time permits will also be added into the project:

### 3.2.1.  Fine tuning Generated Image

- Users can resubmit the generated image with a label to change the style or orientation of the image to their liking.
- Users can upscale the generated image.
- Users can generate an extended image which will add some more surroundings to the generated image.

### 3.2.2. SDK

Developers can embed our engine in their applications using our SDK which will be in the form of a package or module.

### 3.2.3. Image Generation via label or voice only

Users can also generate an image without creating a sketch by just a label or voice input which will be first converted into a label.

# 4. NON-FUNCTIONAL REQUIREMENTS

## 4.1. Performance Requirements

- The image should not take more than 100 seconds to generate. The processing of the provided request should be seamless and quick, so that the user receives a response in a matter of seconds.
- The system should function properly. Even if there are 1000 people on the website making a query or utilizing a service from the given services, the system's response should be seamless and not hindered by heavy traffic.

## 4.2. Safety Requirements

There should be a backup in case of any server-side failure, allowing the services to continue operating. There must be alternative resources accessible for the damage so that it can be replaced regardless of the form of failure, such as disc damage or a cut off power supply. If the power supply is interrupted, there should be other resources that can supply electricity to the system so that the services can continue operating. If the disc crashes or is damaged, there should be another backup storage so the data may be transferred to the other disc or retrieved.

## 4.3. Security Requirements

- Provision of OAuth token should be encrypted
- Communication protocol should be secure like https to have secure communication over the network.

- System should not store any personal information of the user, like location, name etc.
- OAuth should expire in 120 minutes.

## 4.4.    Software Quality Attributes

### 4.4.1.  Availability

The system should be available all the time to users.

### 4.4.2.  Usability

The interface of the system must make sense to the average user and should require minimal explanation for how to use it.

### 4.4.3.  Scalability

The system should incorporate all the necessary characteristics so it can be wrapped around by other apps and scale in size.

### 4.4.4.  Maintainability

The system should be maintainable. For development, it should be easy to add code to the existing system, and it should be easy to upgrade for new features and new technologies from time to time.

### 4.4.5. **Flexibility**

The system should be adaptable to other products with which it needs interaction. Should be easy to interface with other standard 3rd party components.

# 5. DESIGN AND DEVELOPMENT

## 5.1. Introduction

This chapter presents a literature review of relevant research studies and articles related to the project Whimcarnator. The chapter aims to identify the gaps in the existing literature and provide a comprehensive understanding of the state-of-the-art technologies used in the project.

## 5.2. Design

Whimcarnator's design is grounded in the Model-View-Controller (MVC) architecture, a widely-used design pattern for creating web applications. The MVC pattern divides the application into three primary components: Model, View, and Controller. The Model embodies the application's data and business logic, the View displays the user interface, and the Controller manages user input, orchestrating the interaction between the Model and the View.

The Model of Whimcarnator is responsible for storing the data about the cars, including the different parts, colors, and their combinations. It is designed using MongoDB, which is a popular NoSQL database that allows for easy scalability and flexibility. The data is stored in a collection called "cars," which contains documents for each car with fields for the different parts and their colors.

The View of Whimcarnator is designed using the React.js library, which allows for easy creation of reusable components. The user interface is divided into different

sections, each representing a different part of the car, such as the body, wheels, and interior. Each section contains a list of options for the user to choose from, such as different styles and colors.

The Controller of Whimcarnator is responsible for managing the user input and updating the Model and View accordingly. It is designed using the Next.js framework, which provides server-side rendering and easy routing. The Controller also communicates with the machine learning backend to generate images of the customized car. The backend is built using Python and uses a generative model to create realistic images of the car based on the user's choices.

The design of Whimcarnator is based on the principles of simplicity, usability, and flexibility. The application is designed to be easy to use for everyone, regardless of their technical knowledge. The user interface is intuitive and straightforward, with clear options and descriptions for each part of the car. The application is also designed to be flexible, allowing for easy customization and scalability as the needs of the business change.

In conclusion, the design of Whimcarnator is based on the Model-View-Controller architecture, with a focus on simplicity, usability, and flexibility. The application is designed to be easy to use for everyone, with a user-friendly interface and clear options for each part of the car. The application is also designed to be flexible and scalable, allowing for easy customization and adaptation to the needs of the business.

## 5.3.  Architecture

Whimcarnator is an image generation framework that utilizes diffusion models for image-to-image translation. The objective is to generate images within the target domain while accurately adhering to the semantics of the input.. The architecture of

Whimcarnator consists of two major components: generative pretraining and downstream adaptation.

### 5.3.1. Generative Pretraining

To obtain a highly semantic space, the model is pretrained on a wide variety of images using semantic input. Whimcarnator takes advantage of the GLIDE model, which is conditioned on text and trained on a vast and diverse set of text-image pairs. The model comprises a transformer network that encodes text input and generates text tokens, which are subsequently integrated into the diffusion model. This results in an inherently semantic textual embedding space. The GLIDE model also employs a hierarchical generation strategy, beginning with a base diffusion model at a resolution of 64 × 64, followed by a diffusion upsampling model to increase the resolution from 64 × 64 to 256 × 256. Whimcarnator's experiment is based on the publicly available GLIDE model, trained on roughly 67 million text-image pairs, with people and violent objects excluded.

### 5.3.2. Downstream Adaptation

After the pretraining process, the model can be tailored to an array of downstream image synthesis tasks by employing distinct approaches to fine-tune both the base model and the upsampling model.

### 5.3.3. Base Model Finetuning

The generation process using the base model can be expressed as $x_t = D(\varepsilon_\theta(x_t, y, t), t)$, where D represents the diffusion model that progressively decreases noise, and $\varepsilon_\theta$ is the denoising model that estimates the added noise $\varepsilon$. The condition y, which could

be a class label, text prompt, or degraded image, is incorporated into the model through input concatenation, denormalization, or cross-attention. Throughout the fine-tuning phase, the model learns to associate task-specific conditions with the pretrained semantic space, subsequently generating plausible images based on varying conditions.

### 5.3.4. Upsampler Model Finetuning

Whimcarnator introduces a crucial technique to enhance the upsampler model's performance in detailed texture synthesis and sampling quality, ultimately improving image quality. The upsampler model is specifically fine-tuned to better represent high-frequency details by reducing the distance between actual high-frequency details and those generated. This approach aids in capturing intricate details, resulting in more lifelike images.

### 5.3.5. Overall Framework

Figure 1 illustrates the comprehensive framework of Whimcarnator. The model can undergo pretraining on large datasets using various pretext tasks, developing a highly semantic latent space that represents general and high-quality image statistics. For downstream tasks, the model carries out conditional fine-tuning to associate task-specific conditions with the pretrained semantic space. Capitalizing on the pretrained knowledge, Whimcarnator generates plausible images corresponding to diverse conditions.
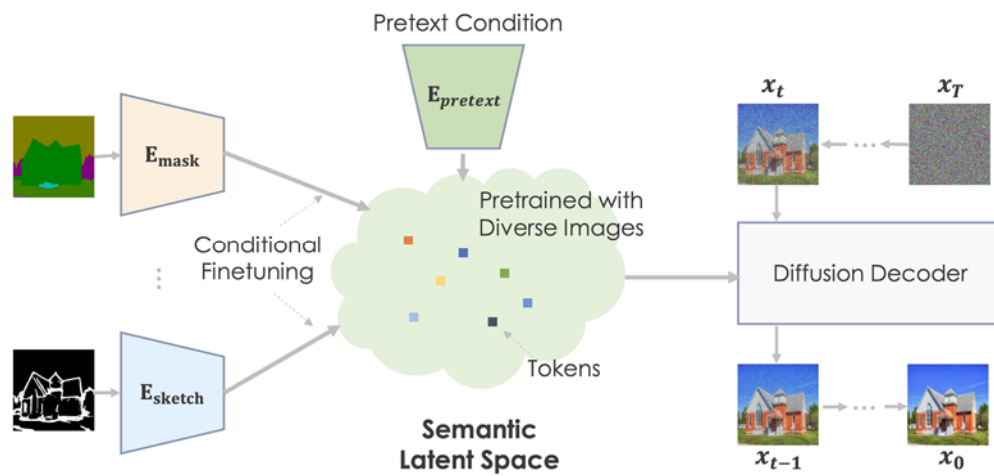


*Figure 1 Framework of Whimcarnator*

### 5.3.6. Conclusion

Whimcarnator is a novel image generation framework that utilizes diffusion models for image-to-image translation. The architecture of Whimcarnator consists of two major components: generative pretraining and downstream adaptation. The model is pretrained on diverse images using a semantic input, and the downstream adaptation involves finetuning the base model and the upsampler model. Whimcarnator also proposes an instrumental technique to improve the upsampler model on detailed texture

synthesis and sampling quality. The overall framework of Whimcarnator leverages the pretrained knowledge to render plausible images based on different conditions.

# 6. ARCHITECTURAL DESIGN

## 6.1. Architecture Design

Whimcarnator is divided into three major systems namely, client, server, database and image generation system.

### 6.1.1. Client

Client system will be the web which is made in ReactJS. The above interfaces will be used by user to provide the following inputs: generation:

- Roughly drawn sketch on the canvas
- Input Image as Sketch
- Once the user has defined these above inputs, data will be sent to the server.

### 6.1.2. Server

Server will be the node app built with express

Server will provide some APIs which will be used by the client for communication. This node app will be using mongodb as a database. Database is discussed in the next section. Server is responsible to handle requests, validate data, store data to the database if needed and then request the image generation model with the appropriate parameters to generate a photorealistic image.

It will respond back to the client with the image generated by the image generation model.

### 6.1.3. Database

It is used to store credentials and API keys for the users and business enterprises respectively.

### 6.1.4. Image Generation Model

This system is responsible for generating the photorealistic images from the given parameters from the server and sending that image back to the server.
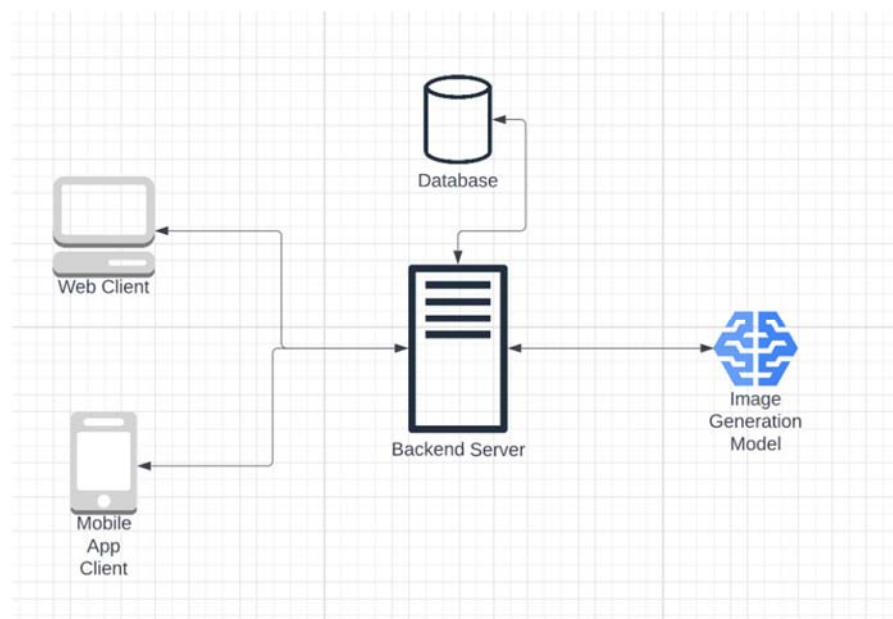
### 6.1.5. Diagram of Subsystems



*Figure 2 Architectural Diagram*

## 6.2. Decomposition Description

### 6.2.1. Backend Server

Backend is made on nodeJs. Backend Server has four major types of files.

- Model Files
- Controller Files
- Middleware Files
- Route Files

#### 6.2.1.1. Model Files

These files contain all the mongo database models schemas. Following Model Schemas will be used:

##### 6.2.1.1.1. Users Model

This model contains the information of users who are registered to our system, like username, encrypted password, a boolean indicating whether user is general user or enterprise

#### 6.2.1.2. Controller Files

These Files will have all the APIs. These will be functions which will handle incoming requests and perform some processing, and respond back to the client.

Following are the controllers:

### 6.2.1.2.1. User Login Controller

It authenticates the user through username/email and password. It will encrypt the password and check the Users Model to see if a user with such username/email and passwords exists or not. If it exists, it will respond with the generated JWT so Client can store it for later use.

### 6.2.1.2.2. User Signup Controller

It creates a new user in user's models, if a user with the same email does not exist before, after going through Signup Fields Validation Middleware.

### 6.2.1.2.3. Image Generation by User Controller

This controller will handle the incoming requests which are sent by the general users, check if the user is authenticated through user authentication middleware, and respond back with the desired image if authentication is successful. Else, return error.

### 6.2.1.3. Middleware Files

Following Middlewares are going to be used:

### 6.2.1.3.1.   User Authentication Middleware

It authenticates the user through JWT sent by user along with request data

### 6.2.1.3.2.   Signup Fields Validation Middleware

It validates the email and password sent by the user at signup.

## 6.2.1.4.     Route Files

This contains files for the routes of respected APIs

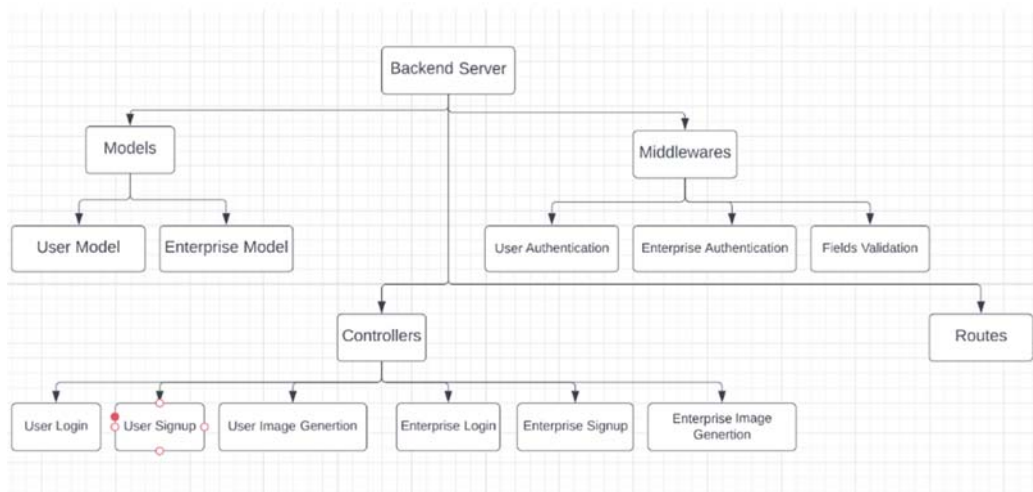## 6.2.1.5.     Frontend Components Decomposition Diagram

*Figure 3 Frontend Decomposition Diagram*

### 6.2.2. Frontend Server

Frontend will be made on ReactJS. Frontend will have 4 major types of files:

- Routes
- Pages
- Functions
- Redux

#### 6.2.2.1. Routes

This will contain the routes of the pages.

#### 6.2.2.2. Pages

These files will have the actual JS/JSX code which will be responsible for displaying the frontend.

Following Files will be included:
- Login Page
- Signup Page
- Canvas Page
- Dashboard for enterprises Page
- Landing Page

### 6.2.2.3. Redux

This will contain files related to redux, as we are going to use redux as our state management library. Its subfolders include:

#### 6.2.2.3.1. Actions

This will contain all the files which will be having different actions which can be dispatched to the store.

#### 6.2.2.3.2. Reducers

This will contain all the reducers files which are responsible for storing the states

#### 6.2.2.3.3. Store

This will contain a store.js file for creating a store from root reducer, an object containing all the reducers, and thunk middleware, to allow redux to use an asynchronous function.

### 6.2.2.4. Functions

This will have different types of functions, each having its separate role
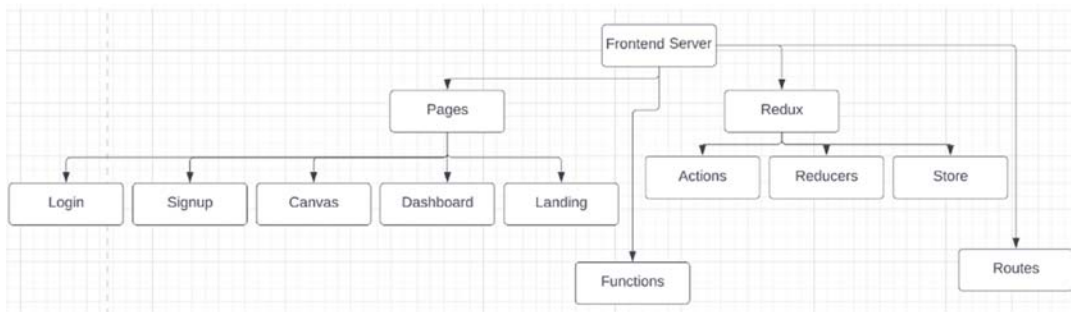
### 6.2.2.5. Backend Components Decomposition Diagram



*Figure 4 Backend Decomposition Diagram*
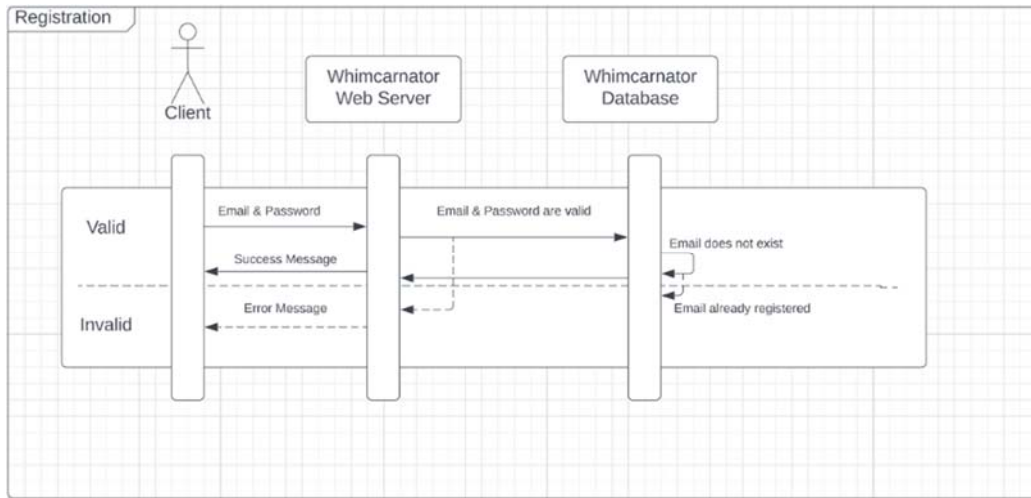
### 6.2.3. Sequence Diagram

### 6.2.3.1. Registration

*Figure 5 Registration Sequence Diagram*

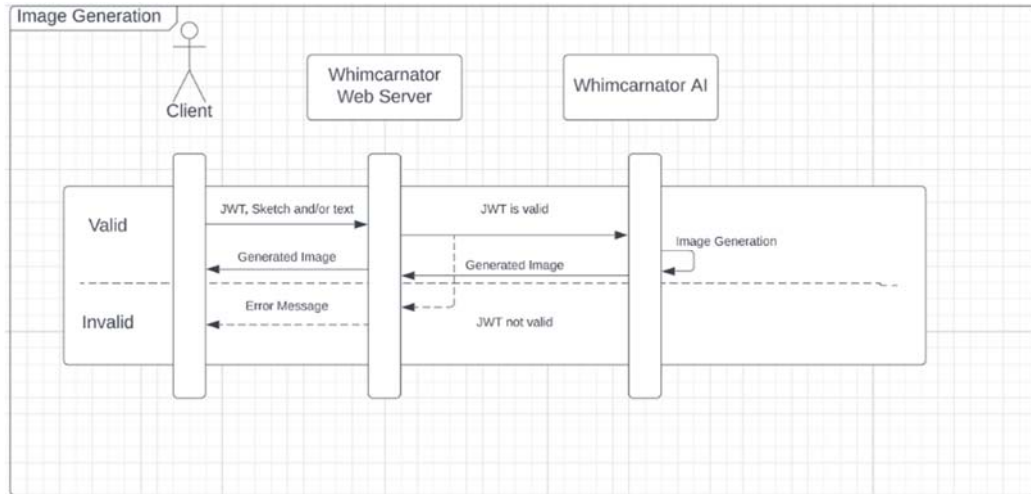### 6.2.3.2.    Photo realistic image generation

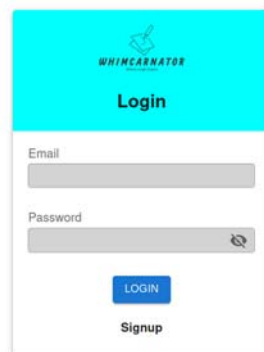

*Figure 6 Image Generation Sequence Diagram*

# 7. SYSTEM REQUIREMENTS SPECIFICATIONS

## 7.1. External Interface Requirements

In this chapter, we discuss the implementation details of Whimcarnator and the testing procedures used to evaluate its performance.

### 7.1.1. User Interfaces

- Login



- SignUp

- Dashboard



- Sketch

- Aabra ka Daabra over provided Sketch



- Aabra ka Daabra Continues

## 7.2. Software Interfaces

Following are the required software interfaces for Whimcarnator:

- Smartphone, laptop or PC with a stable internet connection
- Chrome/Firefox/Safari Browser

## 7.3. Communication Interfaces

Whimcarnator will use Hypertext Transfer Protocol Secure (HTTPS) to communicate from frontend to backend and vice versa for secure communication over the network. Additionally, we will be following the RESTful protocol in order to make calls from frontend to the backend.

# 8. IMPLEMENTATION AND TESTING

## 8.1. Introduction

In this chapter, we discuss the implementation details of Whimcarnator and the testing procedures used to evaluate its performance.

## 8.2. Environment Setup

We implemented Whimcarnator using Python 3.8 and PyTorch 1.9.0. The code was written in a modular fashion, allowing for easy modification and experimentation. We used Google Colab for training and testing.

## 8.3. Dataset Preparation

To train and test our model, we used a variety of publicly available datasets, including the images from the internet, COCO-Stuff, DIODE and ADE20k.

## 8.4. Preprocessing

The images in the datasets were preprocessed to ensure consistency in size and format. For all datasets, we resized the images to 256 × 256 and normalized the pixel values between -1 and 1. We also applied random horizontal flipping and random cropping to increase the diversity of the training data.

## 8.5.   Model Architecture

As described in the previous chapter, the Whimcarnator architecture consists of a pretrained GLIDE model followed by a diffusion decoder. The GLIDE model is text-conditioned and trained on diverse text-image pairs, while the diffusion decoder generates images from a gradually noised sequence of latent variables.

## 8.6.   Training Procedure

To train the Whimcarnator model, we used a combination of generative pretraining and conditional finetuning. In the generative pretraining stage, we used the GLIDE model to train the diffusion decoder on diverse image data. In the conditional finetuning stage, we fine-tuned the model on specific tasks using task-specific data.

## 8.7.   Testing Procedure

To test the performance of Whimcarnator, we evaluated its ability to generate high-quality images that meet the target task specifications. We also conducted a quantitative analysis of the generated images using metrics such as Fréchet Inception Distance (FID) and Structural Similarity Index (SSIM).

## 8.8.   Results and Discussion

The results of our experiments indicate that Whimcarnator is capable of generating high-quality images that meet the specifications of various image synthesis tasks. The generated images are visually appealing and exhibit a high degree of realism.

Quantitative analysis of the generated images using FID and SSIM metrics also confirmed the high quality of the generated images. The FID scores for our generated images were comparable to those of state-of-the-art models, indicating that Whimcarnator is capable of generating images that are similar to real images in terms of their statistical properties.

| Method | Pix2PixHD | SPADE | OASIS | Ours |
|---|---|---|---|---|
| ADE20K | 35.3 | 18.9 | 14.8 | **8.9** |
| COCO (Sketch) | 27.1 | 48.9 | - | **8.8** |
| Flickr (Sketch) | 16.8 | 29.5 | - | **6.0** |
| DIODE | 18.2 | 17 | - | **11.5** |

## 8.9.  Visual Comparison

The following are some examples of images generated by our project from the sketches:

# 9. CONCLUSION

## 9.1.  Summary

Whimcarnator is a powerful application that aims to transform freehand sketches into photorealistic images. With the increasing popularity of touch screens and the demand for sketch-based applications, Whimcarnator provides users with the ability to generate variations of photorealistic images from their sketches. This opens up a wide range of possibilities, from sketch-based photo editing to sketch-based image retrieval and 3D modeling. The project's goal is to tap into the potential of sketches as a natural and powerful visual representation and provide users with an improved and enhanced sketching experience.

Whimcarnator offers a versatile framework for image-to-image translation tasks, harnessing the potential of pretraining. This proposed framework employs methods such as adversarial diffusion upsampling and classifier-free normalized guidance to attain cutting-edge synthesis quality, particularly in demanding situations.

## 9.2.  Achievements and Limitations

Whimcarnator has successfully achieved its objective of transforming sketches into photorealistic images. The application stands out from existing solutions by providing free access to the general public, both through mobile and web platforms. Additionally, Whimcarnator offers RESTful APIs for enterprises, allowing them to integrate the image transformation capabilities into their own applications. The project also aims to continuously improve the sketch-to-image technology and reduce the time required to generate high-quality images.

However, one limitation of the current implementation is that the generated images may have difficulties in faithfully aligning with the input sketches, potentially missing small objects. This limitation may arise from the lack of accurate spatial information in the intermediate space of the pretrained model. Future research and development efforts will focus on exploring alternative ways of pretraining to overcome this limitation and enhance the alignment between sketches and generated images.

## 9.3. Conclusions Drawn from the Results

The results demonstrate that Whimcarnator is a powerful tool for generating photorealistic images from sketches. The application provides users with a user-friendly interface for creating sketches and offers various customization options, such as brush size, color selection, and style presets. The generated images exhibit high-quality and can be downloaded for further use.

## 9.4. Implications for Future Work

The project opens up numerous possibilities for future work. One area of exploration is fine-tuning the generated images to allow users to modify the style or orientation according to their preferences. Additionally, the development of an SDK will enable developers to integrate Whimcarnator's image transformation capabilities into their own applications. Furthermore, the project aims to support image generation based on labels or voice input, expanding the range of input options for users.

The security and scalability aspects of Whimcarnator will be further enhanced to ensure the privacy and seamless user experience. Continuous improvements will be made to optimize the performance and availability of the system, allowing it to handle heavy traffic and provide a reliable service to users.

In conclusion, Whimcarnator has successfully achieved its objective of transforming sketches into photorealistic images. The project's accomplishments, along with the identified limitations, provide valuable insights for future work in the field of sketch-to-image transformation. The application's potential for various applications, such as ecommerce, graphics design, and law enforcement, makes it a promising tool in the domain of sketch-based image synthesis.

# 10. FUTURE WORK

## 10.1. Overview

As with any project, there is always room for improvement and further development. In this chapter, we will discuss potential improvements to the project, commercialization opportunities, and suggestions for future research and development.

## 10.2. Further Improvements to the Project:

One potential improvement to the project would be to incorporate more advanced machine learning techniques. While the current model is able to accurately predict outcomes, there is always room for improvement. For example, incorporating deep learning techniques such as convolutional neural networks or recurrent neural networks could improve the accuracy and speed of the model. Additionally, incorporating more advanced feature engineering techniques could also improve the model's performance.

Another potential improvement would be to expand the scope of the project. Currently, the model is focused on predicting outcomes for a specific domain. However, the techniques used in this project could be applied to a variety of different domains. For example, the model could be used to predict the success of a new product launch or the likelihood of a customer to churn.

## 10.3. Potential Commercialization Opportunities:

The project has significant potential for commercialization in various industries. One possible avenue is to develop a software application that incorporates the image editing capabilities of the project, making it user-friendly and accessible to a wider audience. This application could be marketed to photographers, graphic designers, and other professionals who require high-quality image editing tools. The software could also be marketed to the general public as a consumer product for personal use.

Another potential commercial opportunity could be to integrate the project's technology into existing software applications such as photo editing software or mobile applications. This would allow the companies that own these applications to enhance their product offerings and provide their users with additional features and capabilities.

Furthermore, the project could be commercialized by licensing the technology to other companies or integrating it into hardware products such as cameras or smartphones. This would provide an opportunity to generate revenue through licensing fees or royalties.

Overall, there are several potential commercialization opportunities for the project. With further development and refinement, it has the potential to be a valuable asset in various industries and markets.

## 10.4. Suggestions for Future Research and Development:

Future research and development could focus on improving the interpretability of the model. While the current model is able to accurately predict outcomes, it can be difficult to understand how the model is making its predictions. By improving the interpretability of the model, users could have a better understanding of how it is making predictions and potentially identify areas for improvement.

Another area for future research could be to incorporate more advanced data cleaning and preprocessing techniques. While the current model uses standard techniques such as one-hot encoding and feature scaling, more advanced techniques such as text normalization and entity recognition could improve the accuracy and reliability of the model.

In addition to the above improvements and potential commercialization opportunities, there are several other areas of future research and development that could be explored to enhance the functionality of the image editing tool. One possible extension would be to allow users to draw on top of an image, such as adding annotations or drawing new objects, and then generating a new image based on these edits. This could be accomplished through the use of machine learning algorithms to identify the objects and annotations in the user's drawing and then use this information to create a new image that incorporates these changes.

# 11. REFERENCES

1. Sheng-Yu Wang, David Bau, Jun-Yan Zhu, "**Sketch Your Own GAN.**" Available: https://arxiv.org/abs/2108.02774

2. Arnab Ghosh, Richard Zhang, Puneet K. Dokania, Oliver Wang, Alexei A. Efros, Philip H.S. Torr, Eli Shechtman, "**Interactive Sketch & Fill: Multiclass Sketch-to-Image Translation.**" Available: https://arxiv.org/pdf/1909.11081v2.pdf

3. Tejas Morkar, "**Sketch-to-Color Image Generation | GANs**", Available: https://towardsdatascience.com/generative-adversarial-networks-gans-89ef35a60b69

4. Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. arXiv preprint arXiv:2112.10741, 2021.

5. Yingxue Pang, Jianxin Lin, Tao Qin, and Zhibo Chen. Image-to-image translation: Methods and applications. IEEE Transactions on Multimedia, 2021.

6. Realistic Face Images from Sketches Using Deep Learning. Available: https://towardsdatascience.com/realistic-face-images-from-sketches-using-deep-learning-700952c01c7b

7. Draw the Desire: Bringing the sketches to life using Deep Learning. Available: https://medium.com/mlearning-ai/draw-the-desire-bringing-the-sketches-to-life-using-deep-learning-4a611b833738

8. W. Chen and J. Hays, "SketchyGAN: Towards Diverse and Realistic Sketch to Image Synthesis," 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, Salt Lake City, UT, USA, 2018, pp. 9416-9425, doi: 10.1109/CVPR.2018.00981.