

Conditional Adversarial Image-to-Image Translation with U-Net Generator, PatchGAN Discriminator, and VGG19 Perceptual Loss

**Thesis Submitted
in Partial Fulfillment of the Requirements
for the Degree of**

MASTER OF TECHNOLOGY

**in
Artificial Intelligence**

**by
Tarun Singh
(2K24/AFI/13)**

**Under the Supervision of
Prof. Aruna Bhat**



**To the
Department of Computer Science & Engineering
DELHI TECHNOLOGICAL UNIVERSITY
(Formerly Delhi College of Engineering)
Shahbad Daultapur, Main Bawana Road, Delhi-110042, India**

May 2026



DELHI TECHNOLOGICAL UNIVERSITY
(Formerly Delhi College of Engineering)
Shahbad Daultapur, Main Bawana Road, Delhi-110042, India

CANDIDATE DECLARATION

I Tarun Singh (2K24/AFI/13) hereby certify that the work which is being presented in the thesis entitled "Conditional Adversarial Image-to-Image Translation with U-Net Generator, PatchGAN Discriminator, and VGG19 Perceptual Loss" in partial fulfillment of the requirements for the award of the Degree of Master of Technology submitted in the Department of Computer Science & Engineering, Delhi Technological University in an authentic record of my work carried out during the period from August 2024 to May 2026 under the supervision of Prof. Aruna Bhat.

The matter presented in the thesis has not been submitted by me for the award of any other degree of this or any other Institute.

Tarun Singh

This is to certify that the student has incorporated all the corrections suggested by the examiner in the thesis and that the statement made by the candidate is correct to the best of our knowledge.

Signature of Supervisor



DELHI TECHNOLOGICAL UNIVERSITY
(Formerly Delhi College of Engineering)
Shahbad Daulatpur, Main Bawana Road, Delhi-110042, India

CERTIFICATE BY THE SUPERVISOR

Certified that Tarun Singh (2K24/AFI/13) has carried out their project work presented in this thesis entitled "**Conditional Adversarial Image-to-Image Translation with U-Net Generator, PatchGAN Discriminator, and VGG19 Perceptual Loss**" for the award of **Master of Technology** from the Department of Computer Science & Engineering, Delhi Technological University, Delhi under my supervision. The thesis embodies the results of original work, and studies are carried out by the student himself, and the contents of the thesis do not form the basis for the award of any other degree to the candidate or to anybody else from this or any other University/Institution.

Date:

Prof. Aruna Bhat

Department of Computer Science & Engineering
DELHI TECHNOLOGICAL UNIVERSITY, Delhi, India

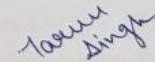
ACKNOWLEDGEMENT

I would like to express my deep appreciation to **Prof. Aruna Bhat** at the Department of Computer Science & Engineering, Delhi Technological University, for her invaluable guidance and unwavering encouragement throughout this research. Her vast knowledge, motivation, expertise, and insightful feedback have been instrumental in every aspect of preparing this research plan.

I am also grateful to **Prof. Anil Singh Parihar**, Head of the Department, for their valuable insights, suggestions, and meticulous evaluation of my research work. Their expertise and scholarly guidance have significantly enhanced the quality of this thesis.

My heartfelt thanks go out to the esteemed faculty members of the Department of Computer Science & Engineering at Delhi Technological University. I extend my gratitude to my colleagues and friends for their unwavering support and encouragement during this challenging journey. I have had some friends that I am thankful to be around. They made me feel truly at home. I would like to thank Mrs. Kiran Bala whom I had such a great time. The last two years were great with such a lovely bunch of people around me. Their intellectual exchanges, constructive critiques, and camaraderie have enriched my research experience and made it truly fulfilling.

While it is impossible to name everyone individually, I want to acknowledge the collective efforts and contributions of all those who have been part of this journey. Their constant love, encouragement, and support have been indispensable in completing this MTech thesis.



Tarun Singh
(2K24/AF1/13)

Conditional Adversarial Image-to-Image Translation with U-Net Generator, PatchGAN Discriminator, and VGG19 Perceptual Loss

Tarun Singh

ABSTRACT

The demand for converting hand drawings into photorealistic images can be attributed to the difficulty of generating rich visuals from sparse, abstract, and incomplete inputs. The rapid growth of creative and design applications by the demand of automation requires better strategies for image synthesis. However, while generative modeling has been proposed as a combination of adversarial training, cycle consistency, and diffusion-based architectures; the use of deep generative systems, improving an architecture time and complexity. Hand drawings are input about a user's coarse sketch. The global movement of having generative models for the public is producing many initiatives. While generative adversarial networks have demonstrated some promising results, there are still challenges, particularly in models trained for conditional generation. Advanced generative techniques in the computer vision domain addressed the critical challenge of preserving semantic layout and ensuring the judicious usage of perceptual losses for models in deep learning. In this model, we have implemented the key techniques which involve Generative Adversarial Networks (GANs), particularly pix2pix, perceptual loss functions, pre-trained VGG-19 network and U-Net architecture. These techniques will provide robust solutions for photorealistic output and secure scene composition. Perceptual metrics are very crucial in providing critical insights into image quality and the mechanics of human-like similarity. We have used a conditional GAN-based architecture (Pix2Pix) with a U-Net generator and perceptual loss, trained on hand-drawn sketches for photorealistic image synthesis. Our results help in demonstrating the effectiveness of the method proposed by offering a scalable solution for the generation of realistic images.

Keywords: Sketch-to-image translation, Generative adversarial networks, VGG-19 network, Image-to-image translation, Perceptual loss, Pix2Pix, Image synthesis.

TABLE OF CONTENTS

Title	Page No.
<i>Acknowledgement</i>	<i>ii</i>
<i>Candidate's Declaration</i>	<i>iii</i>
<i>Certificate</i>	<i>iv</i>
<i>Abstract</i>	<i>v</i>
<i>Table of Contents</i>	<i>vi</i>
<i>List of Tables</i>	<i>viii</i>
<i>List of Figures</i>	<i>ix</i>
<i>List of Symbols and Abbreviations</i>	<i>x</i>
CHAPTER 1: INTRODUCTION	1
1.1 Overview	1
1.2 Problem Statement	2
CHAPTER 2: BACKGROUND AND RELATED CONCEPTS	4
2.1 Generative Adversarial Networks (GANs)	4
2.2 Pix2Pix Framework	5
2.3 U-Net Architecture	6
2.4 Perceptual Loss and VGG-Based Feature Extraction	8
2.5 Image Augmentation Techniques	9
CHAPTER 3: LITERATURE REVIEW	11
3.1 Image-to-Image Translation	11
3.2 Sketch-to-Image Synthesis	12
3.3 GAN-Based Generative Models	13
3.4 Perceptual Loss in Generative Models	14
3.5 Performance Evaluation Metrics	15

CHAPTER 4: PROPOSED ARCHITECTURE	18
4.1 Dataset Description	18
4.2 Data Preparation and Preprocessing	19
4.3 Data Partitioning (Train/Validation Split)	21
4.4 Feature Extraction using VGG19	23
4.5 Model Architecture (U-Net Generator & PatchGAN Discriminator)	24
CHAPTER 5: EXPERIMENTAL EVALUATION	27
5.1 Training Configuration	27
5.2 Validation Process and Testing Process	29
5.3 Results	30
CHAPTER 6: CONCLUSION, LIMITATIONS AND FUTURE SCOPE	33
6.1 Conclusion	33
6.2 Challenges and Limitations	34
6.3 Future Scope	35
CHAPTER 7: SOCIAL IMPACT	36
REFERENCES	38

LIST OF TABLES

Table 2.1. Common image augmentation techniques.

Table 3.1. Performance Evaluation Metrics Used in the model.

Table 4.1. Step-by-Step Pipeline Summary.

Table 4.2. Split Proportions.

Table 4.3. Split Parameters.

Table 4.4. Post-Split Pipeline Behaviour

LIST OF FIGURES

Figure 2.1: Training of a GAN

Figure 2.2. Pix2Pix model

Figure 2.3. U-net architecture.

Figure 2.4. System overview.

Figure 4.1. Feature Extraction Using VGG19.

Figure 4.2. Model Architecture

Figure 5.1. Training Setup.

Figure 5.2. Validation & Testing Setup.

Figure 5.3. Validation & Testing Setup.

Figure 5.4. Training The Model.

Figure 5.5. Final Output.

LIST OF SYMBOLS AND ABBREVIATIONS

GAN	Generative Adversarial Network
G	Generator
D	Discriminator
L1	L1 distance / L1 loss (Manhattan distance)
VGG19	Visual Geometry Group 19
U-Net	U-shaped Network
U-Net++	Enhanced/densely connected variant of U-Net
PatchGAN	Patch-based Generative Adversarial Network
Pix2Pix	Pixel-to-Pixel (image-to-image translation framework)
AR/VR	Augmented Reality / Virtual Reality
ReLU	Rectified Linear Unit
MICCAI	Medical Image Computing and Computer-Assisted Intervention
MSE	Mean Square Error (referenced as "mean square error")
I2I	Image-to-Image (transformation/translation)
cGAN	Conditional Generative Adversarial Network
S2I	Sketch-to-Image
CycleGAN	Cycle-Consistent Generative Adversarial Network
MUNIT	Multimodal Unsupervised Image-to-Image Translation
DRIT	Disentangled Representation Image-to-Image Translation

CHAPTER 1

INTRODUCTION

1.1 Overview

In most machine learning problems, the training data consists of samples of input and output that were collected and labeled beforehand. Such flexibility in obtaining labeled data might not be present when designing algorithms for more creative purposes, since the availability of paired data, like sketches and real-world pictures, can be rare, costly, and specific to the domain. In the case of generating images and applying generative modeling, dedicated algorithms have been designed in order to translate visual information into photorealistic images by means of adversarial training and deep perceptual features.

Moreover, conditional generative models that have been trained using paired datasets of images have shown high effectiveness in image translation applications, including transformation from edge images, semantic segmentation, and sketch to natural images. Given the considerable attention being paid by the computer vision community, human-computer interaction researchers, and digital artists to GANs for visual content generation, a concentrated investigation in this growing field would be highly relevant.

The past few years have seen a number of deep learning techniques applied to image generation, from style transfer to super-resolution to even fully synthetic scenes based on structure information. These breakthroughs have been made possible by large datasets of images, alongside new advances in generative networks, such as U-Net and PatchGAN discriminator. Nevertheless, one issue that has persisted in the task of converting sketches into photographs is that the output should not only be photorealistic in appearance, but also should make sense semantically and perceptually within the domain being targeted. Traditional pixel-based loss functions often result in blurry or inconsistent results that do not match the real-world visual perception.

The task of generating images from sketches is essential in many applications including the creation of digital art, product design, and interactive human-computer

systems. Deep learning methods are especially effective when it comes to identifying important structural information in the form of sparse drawings and translating it into highly dimensional image spaces. The rising need for tools that translate sketches to photographs has resulted in the necessity for reliable generative algorithms that would produce accurate images based on the geometric input of the user's drawing. Yet, creating images of perceptual quality proves to be a difficult task.

The solution proposed by Pix2Pix allows for conditional image-to-image transformations via paired adversarial learning. The method involves the combination of a U-Net network as the generator, which is able to preserve fine details due to skip connections, and a PatchGAN network as the discriminator, which constrains the synthesized samples to look realistic at the local level. Specifically, in our project, we additionally incorporate a perceptual loss term using a pretrained VGG19 network, measuring discrepancy between deep representations as opposed to pixel values. This hybrid loss function, composed of an adversarial, reconstruction, and perceptual component, allows for generating outputs superior to baseline methods.

1.2 Problem Statement

The question of generating realistic images from drawings becomes more complicated by the nature of errors inherent in the traditional methods of image generation. In the case of traditional deep neural network-based image generators, the loss function of the image generator works on the per-pixel basis, e.g., mean square error or L1 distance. Although the method is easy to use mathematically, it ignores perceptual features of a realistic image. The result of the method leads to a famous artifact of blurry images with smoothed contours and textures and semantically incorrect but numerically accurate images.

The problem is aggravated by the vague nature of representation of sketches. It is clear that one hand-drawn sketch can produce several realistic images as sketches do not provide information about color, texture, and illumination. To solve the many-to-one relation problem, the traditional regression approach would simply take an average of the set of generated images, while the outcome does not represent any interpretation of the drawing.

Traditional approach involving generation of images from sketches by use of the

generator network where there is merely pixel level supervision to transform the input sketch into an image cannot possess the capability of imposing higher perceptual consistency. There exist several reasons why images generated by the generator become adversary images characterized by checkerboard effects, color bleeding among others. These arise from the fact that there is no way for the generator to evaluate its output with respect to visual properties of objects in question within the deep learning framework.

There being a rise in demand for automatic approaches to sketch-to-image conversion in various fields including digital art, product design, fashion design and AR/VR applications necessitates a generative model that goes beyond mere pixel matching. Considering the difficulty in training models due to the unavailability of a huge amount of data sets of pairings between sketch and image, it is imperative that much work needs to be done in the domain as it involves manual effort and expertise. Hence the conditional GAN approach, namely Pix2Pix together with VGG19 perceptual loss has been chosen as an efficient approach for sketch-to-image conversion in this research.

CHAPTER 2

BACKGROUND AND RELATED CONCEPTS

2.1 Generative Adversarial Networks (GANs)

Before Goodfellow and colleagues proposed the GAN framework in 2014, training deep generative models was a genuinely difficult problem. Approaches like Restricted Boltzmann Machines and Deep Belief Networks required Markov chain Monte Carlo sampling which is a computationally expensive process prone to poor mixing, while maximum likelihood methods struggled with intractable partition functions. The main issue here was the difficulty of making the neural network not only classify the data, but also learn how to model and generate it based on the original distribution. With GANs, however, the whole task was redefined as a game. The idea behind this model is very simple indeed. Two neural networks are simultaneously trained, one being the generator G and another being the discriminator D . The generator receives random noise as input and tries to produce examples like those belonging to the true data distribution, while the discriminator tries to tell them apart. Goodfellow et al. describe this dynamic using the analogy of counterfeiters and police: the forger gets better at making fake currency, and the detective gets better at catching fakes, and this arms race drives both towards a kind of expertise. Crucially, the entire system is trained end-to-end with backpropagation i.e., no Markov chains, no approximate inference, just gradients.

The training objective formalises this as a minimax game. The discriminator is trained to maximise the log-likelihood of correctly labelling real and fake samples, while the generator is trained to fool it i.e., minimising $\log(1 - D(G(z)))$ or equivalently maximising $\log D(G(z))$ to avoid gradient saturation early in training. Figure 1 below illustrates how this plays out over the course of training [2].

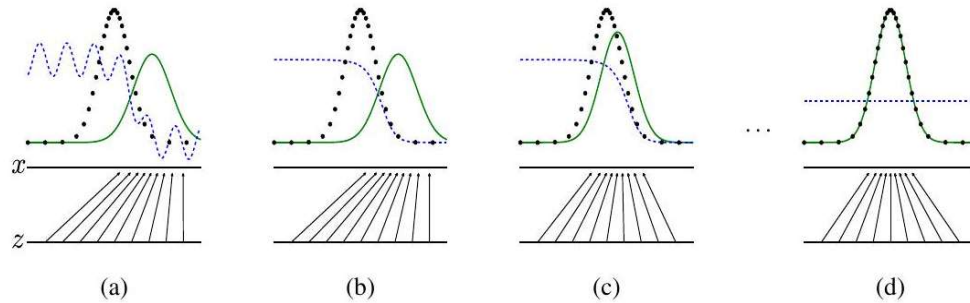


Figure 2.1: Training of a GAN. Adapted from [1].

2.2 Pix2Pix Framework

Image-to-image translation has been one of the most active and practically useful areas of computer vision for the past decade. The basic idea is straightforward: given an image in one domain, say a pencil sketch it learns a mapping that produces the corresponding image in another domain, such as a realistic photograph. For much of this history, such models were task-specific, requiring bespoke architectures trained on narrow datasets. The introduction of Pix2Pix by Isola et al. in 2017 [22] changed that by providing a general-purpose framework built on conditional GANs that could handle almost any paired image translation task with a single architecture. The original Pix2Pix uses a U-Net as its generator, an encoder-decoder network with skip connections that help preserve spatial information during the encoding and decoding process. A PatchGAN discriminator — rather than judging the whole image at once — evaluates overlapping patches of the output, which encourages local realism and texture consistency. The training objective combines a conditional GAN loss with an L1 pixel-wise penalty to keep generated images close to their ground-truth targets. Despite its strengths, the model has two notable failure modes. First, U-Net's skip connections are constrained to link only layers of the same scale, which limits how effectively multi-scale features can be combined and leads to some unnecessary information loss. Second, the standard discriminator only sees the full generated image and its source as it has no specific mechanism to penalise the model for getting wrong precisely the parts that changed between the source and target domain. The paper [22] proposes two targeted improvements to address these issues. The first is to swap the

U-Net generator for a U-Net++, a more densely connected variant introduced by Zhou et al. in 2020 that allows skip connections between feature maps at different scales and depths — not just the same level. This means the decoder can draw on a richer mixture of coarse and fine features simultaneously, reducing information loss and producing crisper, more detailed outputs. In addition to the architectural change, the authors modify the internal convolutional blocks of the U-Net++ (replacing VGG-style classification blocks with strided convolution, channel normalisation, and dropout layers) to make it better suited to the image translation setting. The second improvement is more conceptually novel: the introduction of a dedicated differential image discriminator. As shown in Figure 2 [22] below, the standard Pix2Pix already includes a generator and a PatchGAN discriminator; the authors add a second discriminator that operates not on the generated image itself, but on the difference image — the pixel-wise subtraction of the source input from the target.

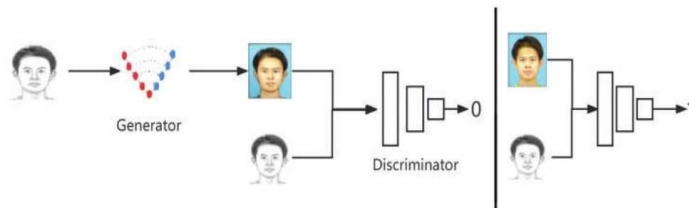


Figure 2.2. Pix2Pix model. Adapted from [22].

2.3 U-Net Architecture

Applying deep learning to image segmentation presents a very different challenge to standard computer vision tasks. In classification, a single label is assigned to an entire image; in segmentation, every pixel must be labelled, and in medical imaging that means the model needs to understand both global context — what structure it is looking at — and fine local detail simultaneously. The other major obstacle is data scarcity. While ImageNet offered millions of labelled training examples to the computer vision community, biomedical researchers typically have access to only tens or hundreds of annotated images. Before U-Net, the leading approach was a sliding-window convolutional network that made pixel-wise predictions by evaluating small patches around each pixel independently — a method that was accurate but painfully slow, redundant, and caught in an irresolvable tension between seeing enough context and preserving localisation detail. Ronneberger, Fischer, and Brox set out in 2015 to

break that tension entirely.

Their solution, which they presented at MICCAI 2015, extended the concept of a fully convolutional network into what they called a u-shaped architecture — the name U-Net comes from the distinctive shape it traces in a diagram. The key insight is that the network is split into two symmetric halves: a contracting path on the left that progressively compresses the input through repeated 3×3 convolutions, ReLU activations, and 2×2 max-pooling steps, doubling the number of feature channels at each stage while halving the spatial resolution; and an expansive path on the right that reverses this process, using transposed convolutions to upsample the feature maps back to the original resolution while halving the channel count at each step. What makes the architecture work is that the two halves are connected by skip connections — at each depth level, the feature maps from the contracting path are copied and concatenated directly into the corresponding level of the expanding path. This means the decoder always has access to high-resolution spatial information from early in the network, not just the compressed bottleneck representation, allowing it to reconstruct fine boundary detail that pure encoder-decoder networks tend to lose.

Figure 3 below illustrates this architecture clearly. The full network has 23 convolutional layers in total, no fully connected layers at all, and ends with a 1×1 convolution that maps the feature vectors at each pixel position to the desired number of output classes. Because the convolutions are unpadded, the output segmentation map is slightly smaller than the input, which is why an overlap-tile strategy is used for large images — neighbouring tiles are processed with mirrored padding at their borders so that no context is lost at the edges [3].

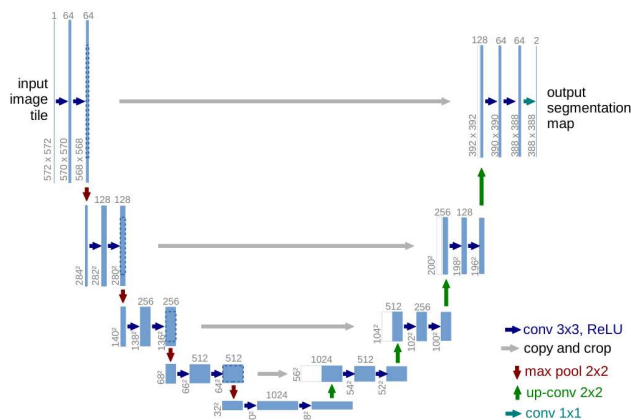


Figure 2.3. U-net architecture. Adapted from [3].

2.4 Perceptual Loss and VGG-Based Feature Extraction

Johnson et al. (2016) tackle a fundamental challenge in image processing: how to make neural networks that transform images both fast and visually convincing. Traditional approaches to tasks like artistic style transfer and image super-resolution tend to fall into one of two camps. On one side, feed-forward convolutional neural networks can process images in real time, but they rely on per-pixel loss functions that measure quality by comparing individual pixels between the output and a reference image. This turns out to be a poor proxy for how humans actually perceive image quality — two images that look nearly identical to the eye can differ dramatically pixel-by-pixel if one is slightly shifted. On the other side, optimization-based methods like that of Gatys et al. produce stunning results by measuring image similarity through high-level features extracted from a pretrained network, but doing so requires running hundreds of iterations of gradient descent for every single image, making them far too slow for practical use. The key insight of this paper is that these two approaches do not have to be in conflict. The authors propose training a feed-forward transformation network using perceptual loss functions — losses defined not over raw pixels, but over feature activations drawn from a fixed, pretrained VGG-16 network. Because VGG-16 was trained for image classification, its internal representations already encode rich semantic and perceptual information about content and style. By asking the transformation network to produce outputs whose VGG features resemble those of a

target image, the authors effectively transfer that semantic knowledge into a fast, single-pass network. For style transfer, this approach yields results that are qualitatively comparable to Gatys et al. while being roughly a thousand times faster, enabling real-time stylization at 20 frames per second. For single-image super-resolution at $\times 4$ and $\times 8$ scale factors, replacing the standard per-pixel loss with a perceptual loss leads to sharper edges and finer recovered details — particularly in semantically meaningful regions like faces and textures — even though standard metrics like PSNR and SSIM score it lower, highlighting a known disconnect between those metrics and human visual judgment. Together, these results demonstrate that perceptual loss functions are a powerful and general training signal for image transformation tasks, bridging the gap between speed and perceptual quality [4].

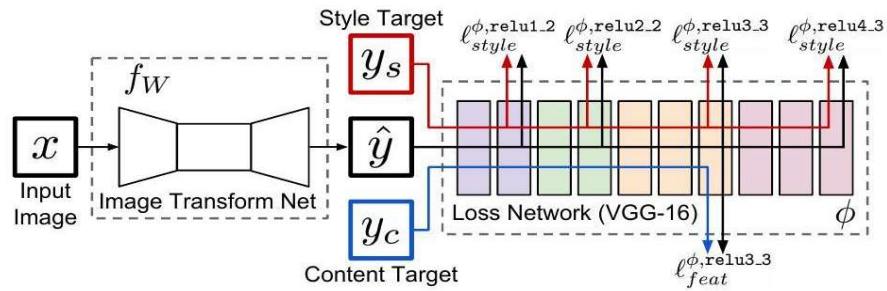


Figure 2.4. System overview. Adapted from [4].

2.5 Image Augmentation Techniques

Image augmentation involves the use of several methods that could be used within the context of machine learning or computer vision in order to artificially increase the size and diversity of a specific training dataset using a number of manipulations that could be conducted on images. Using such a method offers an alternative solution to how a model will be able to experience a variety of scenarios in terms of images, without having to gather any extra data, since this process can be costly and time-consuming. Typical manipulations associated with image augmentation would involve geometric transformations (for example, flipping, rotating, cropping and scaling), photometric transformations (such as changing brightness, contrast, saturation, hue), as well as some more advanced manipulations, such as adding Gaussian noise, elastic transformation, random erasing of parts of an image or combining two different

images, among others, using various approaches (such as Mixup and Cutmix). The basic idea behind this method is to prevent the risk of overfitting through avoiding dependence on specific pixel patterns, while increasing the ability to generalize.

Table 2.1. Common image augmentation techniques

Technique	Description	Use Case	Effect on Dataset
Horizontal Flip	Mirrors image left-to-right	Object detection, classification	Doubles data size
Random Rotation	Rotates image by random angle	Medical imaging, aerial views	Increases rotational variance
Zoom / Cropping	Zooms in or crops a region	Face recognition, fine-grained tasks	Emphasizes local features
Color Jitter	Alters brightness, contrast, saturation	Outdoor scene recognition	Improves lighting robustness
Gaussian Noise	Adds random pixel-level noise	Sensor noise simulation	Improves noise tolerance
Cutout / Erasing	Randomly masks rectangular regions	Occlusion robustness tasks	Forces feature diversity
MixUp	Blends two images with weighted average	Classification, regression	Encourages smoother predictions
Elastic Distortion	Applies random elastic deformations	Handwriting, medical scans	Models morphological variation

CHAPTER 3

LITERATURE REVIEW

3.1 Image-to-Image Translation

Architecture of I2I transformation processes has been developed into a paradigm for revolutionary improvements in computer vision through the integration of numerous pixel manipulation approaches into one generative model[5]. The Pix2Pix framework was the first supervised architecture to utilize cGAN and an element-wise L1 loss function to ensure precise low-frequency structure alignment in conjunction with PatchGAN-based discriminators for capturing sharp high-frequency structures [5, 7]. In order to overcome the limitations associated with auto-encoders, particularly the information bottleneck, the paradigm introduced the idea of using U-Net structure for a generator, which is based on the idea of applying spatial skip connections between parallel layers allowing direct passage of critical feature maps [13, 21]. In consideration of the extremely complex task of collecting large-scale pixel-matching image sets, unsupervised architectures became increasingly common following the development of CycleGAN[6]. The unsupervised CycleGAN managed to establish a mapping between two independent, unordered sets of domains (X to Y and Y to X) through a mathematical cycle consistency criterion where $F(G(x))$ approximately equals x , effectively constraining the search space for solutions and maintaining invertibility of structural anatomy without direct pairing [6, 15].

Nevertheless, bidirectional cycle-consistency constraints being deterministic in nature inevitably lead to mode collapse, that is, a one-to-one mapping between an individual source sample and an equally fixed target sample [8]. Multimodality, which is highly required in practice, has led to the advent of sophisticated networks for feature disentanglement such as MUNIT and DRIT, where the latent space is split into content space shared by domains corresponding to pose and attribute space domain-dependent on style variations like lighting and texture [8, 9]. Contemporary

approaches have gone a step further in response to drastic distribution discrepancies and spatial deformations, such as converting abstract sketches to realistic photographs and cross-domain synthesis in the field of medical imaging, such as from MRI to CT maps, through the use of more reliable self-supervised semantic bridges and diffusion bridge frameworks rather than instable adversarial losses [18, 19]. The contemporary systems rely heavily on pre-trained visual encoders that extract robust geometric information regardless of appearance [19, 21].

3.2 Sketch-to-Image Synthesis

The history of sketch-to-image (S2I) architectural progression can be considered a highly specific and technically complicated sub-domain within an image-to-image translation branch due to converting highly abstract black-and-white drawings to photorealistic textures and background images[21]. The area started developing using a Pix2Pix conditional generative framework, applying a conditional GAN model along with a combined loss of L1 reconstruction on pixel-level and adversarial loss to cover a large gap between binary sketches and color photos [5, 10]. In order not to destroy the sharp outline boundaries during learning, initial architectures used a layered U-Net architecture, connecting features in parallel feature maps in a skip connection layer, thus avoiding compression limitations inherent to conventional bottleneck auto-encoders [16, 21]. However, since creating pairs of aligned images in the form of sketch and picture requires significant human resources either manually or automatically (involving filters for edge detection but without the human style of drawing), it soon became obvious that unsupervised approaches to mapping images needed to be developed [6, 21]. The base of these models was CycleGAN, applying two unpaired collections of drawings mapping them to the photo domain (X to Y and Y to X) with the help of a cycle-consistency loss function where $F(G(x))$ approximately equals x [6, 9]. This technique proved to be very effective in penalizing structural drift and forced the network to respect the initial contour lines regardless of not having explicit target pairs [6, 15].

While demonstrating significant improvements in structural consistency, standard cycle-consistent methods still worked deterministically and suffered greatly from a

catastrophic mode collapse problem, producing just one deterministic photo translation from each input sketch [8]. Given that a simple contour sketch can be associated with an infinite number of possible texture/material combinations and lighting conditions, scientists had to design multimodal feature disentanglement networks, such as MUNIT and DRIT [8, 9]. This approach entailed the separation of latent space into a content part (preserving the spatial configuration of sketch lines) and an attributes part (accounting for style variations, colors, and shades), thus allowing to generate various images based on the same drawing configuration [8]. As artificial intelligence entered the modern era, classical GAN-based architectures started receiving criticism for unstable training procedures, structural deformations, and blurriness of the output in case of complicated and non-rigid inputs [18, 19]. Nowadays, advanced S2I architectures solve this issue by using pre-trained self-supervised semantic bridge encoders together with latent diffusion bridge architectures [15, 30]. In utilizing powerful feature descriptors (like DINO Patch Embedding) that disregard the effect of surface texture changes to capture essential geometries, these contemporary Diffusion Models have the ability to gradually denoise towards hyper-realistic, text-driven images while remaining perfectly structurally consistent with the initial hand drawings provided by users [21].

3.3 GAN-Based Generative Models

Structural transformations of GAN-based similarity evaluation models have been a breakthrough in image-to-image translation, shifting focus from simple pixel-to-pixel loss calculations to more sophisticated structural alignment methods based on deep architecture [1]. For example, traditional pixel-wise loss such as L1 and L2 could not adapt to domain changes because their objective was to compute the absolute differences between individual pixels, resulting in blurriness and low-frequency filtering artifacts [5, 7]. To address that problem, pioneering models introduced a deep perceptual similarity measure, using a pre-trained classification backbone, such as VGG-16, as a loss function for generating more realistic images [4, 10]. Specifically, by comparing the activations of corresponding convolutional layers of the generated and reference images, it became possible to maintain

structural integrity, rather than aligning actual pixels [1].

Furthermore, the PatchGAN concept, which serves as a structural similarity classifier based on evaluating the realness of patches in an image matrix, ensured that the generator optimized for sharp and high-frequency textural coherence within localized image patches [5, 17]. With the transition to unsupervised and unaligned image spaces, it has become impossible to solely rely on one-sided perceptual losses due to the lack of target-pair relationships, resulting in extreme geometric distortions [5, 15]. To ensure stable unaligned distribution mappings, cycle-consistency concepts in architectures such as CycleGANs and stacked cycles were adopted, where the network needed to map an input image to another domain ($G(x)$) and then back to its original state from the generated image ($F(G(x))$ approximately x), effectively acting as a global structural similarity constraint [6, 20]. Indeed, the objective function of cycle consistency constrains the scope of solutions to a significant degree because of the stringent criteria set forth [15]. To balance the interplay between the inputs and outputs from the two domains, modern GAN similarity networks resort to the patch-wise contrastive learning framework [14]. The mutual information across the selected layers is maximized based on certain relations of patches as well as through query-selective attention maps that enable these systems to match patterns from a new domain in an incredibly precise manner.

3.4 Perceptual Loss in Generative Models

The introduction of perceptual loss functions into structural and style similarities was an important achievement in the development of generative models because the latter method was much more advanced than the former, as it took more than just comparing pixel matrices [4, 1]. The traditional methods utilized simple mathematical equations to represent the objective of the training process, including L_1 absolute error or L_2 squared error [5]. Even though these objectives ensure global consistency, the loss measurement is based on comparison of only those pixel matrices that are located at identical coordinates, making the slightest shift very impactful [4, 7].

This problem led to introducing the concept of perceptual loss, wherein the hidden

activations are extracted from pre-trained networks, such as VGG-16 and VGG-19, to evaluate the image according to its high-level characteristics [10]. Instead of determining how similar the values of pixels are, the loss function uses two input images: one created using the model, and another – the ground truth [4]. Thus, when calculating the Euclidean distances between the feature maps of two images in various convolutional layers, it ensures structure and semantic accuracy rather than pixel similarity [1].

But the above approach is further subdivided into two main losses, that is, feature reconstruction loss which preserves the overall structure and shapes of the images, and style reconstruction loss which uses the Gram matrix for identifying and reconstructing features such as colors and textures [7]. When it comes to contemporary architectures, the above theory has led to the development of IQA methods such as LPIPS and DISTs [1]. In terms of the application of these metrics, it requires using deep convolutional networks to evaluate the diversity of texture and structure accuracy of the patches [1]. However, by combining these losses with PatchGANs discriminators focused on individual patches, it has enabled today's generator architecture to achieve a good balance [17].

3.5 Performance Evaluation Metrics

The creation and development of generative adversarial models rely on the existence of a good ecosystem of loss functions, which are responsible for the harmony between the generator and discriminator [5]. In this case, the main component of such optimization is Discriminator Loss, which consists of the mathematical expression of the model's capability to optimize its scoring on real images and differentiate synthesized images from real ones [5]. In turn, GAN Loss (adversarial loss function) shows how the generator attempts to minimize it through deception [5, 7].

Because pure adversarial training is highly fluid and prone to spatial distortions, stable image-to-image pipelines add explicit similarity regularizers to anchor geometric reconstruction [6, 15]. The first of these, L1 Loss, calculates the absolute, pixel-by-pixel difference between the generated image matrix and the ground-truth

target, enforcing low-frequency structural alignment and preventing global content drift [7,]. To prevent the blurry artifacts that happen when relying entirely on pixel-level smoothing, frameworks integrate Perceptual Loss [4]. This metric extracts activation maps from fixed, hidden convolutional layers of a pre-trained network (like VGG-16) to minimize high-level semantic and textural differences rather than literal coordinate offsets [4].

These specific objectives are then aggregated mathematically utilizing weighting factors to formulate the overall objective referred to as Total Gen. Loss (Total Generator Loss) used to optimize our model [4, 5].

In terms of monitoring the training stability as well as the ability to generalize throughout the entire process, the overall objective will be tracked in the two different processes: Train Gen. Loss and Val Gen. Loss. While the former tracks the direct performance in terms of convergence and optimization in the current training batch [5], the latter computes the same total loss equation on an entirely different, unseen partition [17]. Monitoring the difference between the two can ensure the prevention of overfitting, making sure that the model learns to find the right cross-domain mappings rather than just memorizing the input dataset textures [16, 17].

Table 3.1. Performance Evaluation Metrics Used in the model.

Metric	Component	Formula / Method	Purpose
GAN Loss	Generator	Binary Cross-Entropy vs. ones	Fools discriminator
L1 Loss	Generator	Mean Absolute Error ($\times 100$)	Pixel-level accuracy
Perceptual Loss	Generator	VGG19 feature MSE ($\times 10$)	High-level realism
Total Gen. Loss	Generator	GAN + $100 \times L1$	Overall

Metric	Component	Formula / Method	Purpose
		+ 10×Perceptual	generator quality
Discriminator Loss	Discriminator	Real loss + Generated loss	Distinguishes real vs. fake
Train Gen. Loss	Training	Avg over batches / epoch	Tracks training progress
Val Gen. Loss	Validation	Avg over val set / epoch	Detects overfitting

CHAPTER 4

PROPOSED ARCHITECTURE

4.1 Dataset Description

At the core of the study is the CUFS database (CUHK Face Sketch Database), an elaborately curated repository that has earned itself a name as one of the most reliable standard databases for face sketch synthesis and recognition studies. Composed of 606 pairs of face images, which consist of images sourced from three different data sets; namely, 188 from CUHK student faces database, 123 from AR face database and 295 from XM2VTS database, what makes this particular repository unique is not the quantity, but rather the pairing, as an experienced artist drew the sketch on a pencil paper based on each photograph he saw, all of them captured in a frontal view with normal lighting and facial expressions. Such an artificial creation by humans makes the images more realistic than any purely algorithmic method. Concerning the current study, the CUHK database will be applied to train a Pix2Pix generator with 80% of paired data used for training and 20% for validating the model. The images will be rescaled to 256×256 pixels and normalized to pixel values from minus one to one in order to ensure stable gradients while training. Live augmentation (flipping the images horizontally and adjusting their lightness randomly) helps enrich the dataset with visual variations without needing to collect any additional images. Thus, the task assigned to the network involves learning the correlation between the artistic style of drawing and real-life photographs since this correlation is essential to transform the sketch into a real photo and requires deep understanding of how the human face should look in accordance with the particular artistic features.

Total Subjects: 188 paired sketch-photo samples

Data Sources: CUHK Database

Sketch Method: Hand-drawn by a trained artist while viewing the photo

Image Conditions: Frontal pose · neutral expression · normal lighting

Image Resolution: 256×256 pixels (resized)

Normalization: Pixel values scaled to $[-1, 1]$

4.2 Data Preparation and Preprocessing

Before a single weight in the network could be updated, the raw CUHK sketch-photo pairs had to pass through a carefully designed preparation pipeline — one built to be fast, reproducible, and kind to a GPU that has better things to do than wait for disk reads. The process begins at ingestion: the user uploads a ZIP archive through Google Colab's file interface, which the code unpacks in place, landing two sorted folders — one for sketches, one for photographs — in the runtime's working directory. From there, the pipeline checks that every sketch has an exact partner, splits the pairs into an 80/20 training-validation partition using a fixed random seed for reproducibility, and hands each file path to a `load_image` function that decodes the JPEG, resizes it to 256×256 pixels, converts any stray grayscale image to three-channel RGB, and scales every pixel from the standard 0–255 range down to the $[-1, 1]$ interval that GAN training strongly prefers. On the training side, augmentation adds one more layer of variation — each batch sees randomly flipped images and subtly shifted brightness values, applied identically to the sketch and its target photo so the pairing is never broken. The finished datasets are batched at 16, shuffled during training, and prefetched using TensorFlow's AUTOTUNE so the GPU stays fed rather than idle between steps.

Table 4.1. Step-by-Step Pipeline Summary

Step	Operation	Parameter	Purpose
1. Data Ingestion	ZIP file uploaded via Google Colab; extracted to <code>./sketches</code> and <code>./images</code> folders	—	Organise raw paired data on runtime
2. File Pairing	Paths sorted alphabetically; paired positionally; assert	—	Guarantee sketch to photo correspondence

Step	Operation	Parameter	Purpose
	statement verifies equal count		
3. Train/Val Split	sklearn train_test_split with stratified random shuffling	test_size=0.2, random_state=42	Reproducible 80/20 partition
4. JPEG Decoding	tf.io.read_file → tf.image.decode_jpeg to tensor	—	Load raw pixel data into memory
5. Resizing	tf.image.resize applied to every image	256 × 256 px	Uniform spatial dimensions for batching
6. Grayscale → RGB	tf.image.grayscale_to_rgb when channel count == 1	—	Ensure 3-channel input for VGG19 & GAN
7. Normalisation	(pixel / 127.5) - 1 applied per channel	[-1, 1]	Centre activations; stabilise GAN training
8. Augmentation	Random horizontal flip + random brightness on both sketch and photo simultaneously	max_delta = 0.2	Reduce overfitting; add lighting variance

Step	Operation	Parameter	Purpose
9. Batching & Shuffle	Training set shuffled (buffer=100) then batched; validation set batched only	batch = 16	Stochastic mini-batch gradient descent
10. Prefetching	tf.data AUTOTUNE prefetch on both train and val datasets	AUTOTUNE	Overlap I/O with GPU compute

4.3 Data Partitioning

After uploading and extracting the CUHK dataset, the next step is determining which files will serve as the basis of learning and those used to evaluate model performance which is a seemingly trivial yet highly impactful decisions. Instead of separating these datasets manually, the pipeline uses the scikit-learn `train_test_split` method to shuffle the complete list of 606 pairs of file paths and allocate 20% for validation and 80% for training purposes. Importantly, the division affects not only the sketch paths but also the image file paths, ensuring that no sketch and its respective photo get separated and each sketch always goes with the corresponding photo. The random state equals 42, indicating that the separation is the same across all executions of the notebook. No test data set is created in the code; however, this validation split plays dual roles, as it serves both as a means of early stopping during training and helps with comparing the final visual outputs.

Table 4.2. Split Proportions

Subset	Proportion	Approx. Samples (606 total)
Training Set	80%	~485 paired

Subset	Proportion	Approx. Samples (606 total)
		samples
Validation Set	20%	~121 paired samples
Test Set	None defined	—

Table 4.3. Split Parameters

Parameter	Value & Rationale
Function used	sklearn.model_selection.train_test_split
test_size	0.2 — reserves 20% of pairs for validation
random_state	42 — fixes the shuffle seed for reproducibility
Shuffle	Implicit (default True) — randomises order before splitting
Stratification	Not applied — no class labels in this paired regression task
What is split	File path lists (sketch_paths, image_paths) — not raw pixels
Pairing preserved	Yes — sketch and image lists split with identical indices

Table 4.4. Post-Split Pipeline Behaviour

Dataset	Shuffle	Augmentation
train_dataset	Yes (buffer = 100)	Defined but not applied*
val_dataset	No	None

4.4 Feature Extraction using VGG19



Figure 4.1. Feature Extraction Using VGG19.

At the heart of this project lies a very straightforward idea, instead of teaching a neural network everything from scratch, we use the already available one. VGG19 is a deep convolutional neural network originally trained on over a million images from ImageNet which already knows how to see the world. It understands edges, textures, shapes, and high-level visual concepts without us having to show it a single example from our own dataset. So rather than reinventing the wheel, we load VGG19 with its pretrained weights, freeze it completely (no learning, no updates), and tap into its `block5_conv2` layer where a rich, abstract level where the network has already distilled raw pixels into meaningful patterns. This becomes our dedicated Feature Extractor.

When we feed both a generated image and the real target image through this extractor, it returns two sets of feature maps — essentially two fingerprints of visual content. We then measure how far apart those fingerprints are using mean squared error, producing what is known as the Perceptual Loss. This loss doesn't care about pixel-level perfection — it cares about whether the generated image feels right: whether its structure and semantics match the ground truth at a deeper level. Inside our Pix2Pix GAN, this perceptual loss works in tandem with the adversarial GAN loss and a pixel-level L1 loss, forming a three-part training signal that pushes the generator to produce images that are not only realistic to a discriminator, but also structurally faithful and visually coherent. Input sketches and target images are first resized to 256×256 pixels, normalized, and augmented with random flips and brightness shifts before being fed into the pipeline. The generator (a U-Net architecture) then attempts to bridge the gap between sketch and reality, while VGG19 quietly acts as the visual conscience — ensuring every generated image carries the right visual DNA. The result is a model that can transform rough hand-drawn sketches into photorealistic images.

4.5 Model Architecture

In summary, the proposed model is a Generative Adversarial Network based on the pix2pix framework and focuses on translating the hand-sketch images into realistic ones. Differently from CycleGAN, which processes the unpaired images and has to

discover their correlations on its own, the proposed model adopts a more straightforward strategy i.e., learns on paired images, where the sketch image and real one are provided together. With the use of perceptual losses, the training is simplified for the generator producing more accurate images.

Generator: A U-Net architecture serves as a backbone of the generator, which basically implies using an encoder-decoder network with skip connections between matched layers of each part of it. Such design helps preserve the finest spatial features by transferring them between the two parts of the network via skip connections, instead of compressing them at the bottleneck. The generator makes use of instance normalization in the course of processing a 256×256 RGB sketch into a 256×256 RGB image.

Discriminator: Rather than evaluating the realism of the entire picture, the discriminator follows the architecture of PatchGAN, focusing on local patches. The output feature map size is 30×30 , with values indicating the realism of each of them separately, compelling the generator to generate a picture with high-quality details instead of good composition as a whole.

Loss Function: Three different loss functions used for training complement each other. An adversarial loss encourages the generator to generate pictures good enough to pass the discriminator test, trained with binary cross-entropy. At the same time, an L1 loss (multiplied by 100 times), forces the model to pay attention to the big structure, rewarding it when the differences between images on a per-pixel basis are small. Finally, the perceptual loss compares the output of the generator and target image features based on the representation obtained with pre-trained VGG19 model through block5_conv2, comparing their representations via mean squared error criterion.

Training

Setup

Both the generator and discriminator networks have been trained using Adam optimization, where the learning rate is set to $2e-4$ and β_1 is set to 0.5, both of which make an ideal combination for GANs. Batches of 16 images have been used to train the network, and data augmentation has been done using flipping and changing image brightness levels.



Figure 4.2. Model Architecture.

CHAPTER 5

EXPERIMENTAL EVALUATION

5.1 Training Configuration

The generator and discriminator are trained using the Adam optimizer, which has proven successful for GAN-based architectures due to the capability to deal with noisy gradients and the possibility to adapt learning rates dynamically. The learning rate is defined as $2e-4$ with a value of $\beta_1 = 0.5$ (unlike in the Adam optimizer, for which the default value is 0.9). A smaller value of β_1 makes it easier to avoid overshooting and oscillations, which would happen if the model learned too much of previous gradients and became too eager to change the weights accordingly.

A batch size of 16 is used during training; it can be considered an acceptable compromise between the need for memory optimization and gradient stabilization. In order to increase generalization capabilities of the model, each pair of images undergoes a series of augmentations before being used as input data. The list includes random flipping (with a horizontal orientation only) and random brightness adjustments (with the max_delta parameter being set at 0.2). Both an input sketch and its target image are being augmented at the same time in order to preserve their relationship.

Training and validation datasets are split in the ratio of 80 to 20 using the random split function from TensorFlow library with a seed value of 42. All the images, both sketch and targets, are resized to 256×256 pixels and are rescaled to have pixel values in the range from -1 to 1 by dividing all the pixel values by 127.5 and subtracting one. Any grayscale image gets automatically transformed to RGB in order to make sure that the images have three channels as required by the network architecture. The training pipeline is constructed using `tf.data` module of TensorFlow with parallelization of loading data, shuffling of data, batching, and prefetching in order not to slow down the GPU.

```

train_dataset, val_dataset = load_dataset(sketch_folder, image_folder)

generator = pix2pix.unet_generator(3, norm_type='instancenorm')
discriminator = pix2pix.discriminator(norm_type='instancenorm')

loss_object = tf.keras.losses.BinaryCrossentropy(from_logits=True)

def perceptual_loss(target, prediction, feature_extractor):

    target_features = feature_extractor(target)
    prediction_features = feature_extractor(prediction)

    return tf.reduce_mean(tf.square(target_features - prediction_features))

def generator_loss(disc_generated_output, gen_output, target):
    gan_loss = loss_object(tf.ones_like(disc_generated_output), disc_generated_output)
    l1_loss = tf.reduce_mean(tf.abs(target - gen_output))
    perceptual = perceptual_loss(target, gen_output, feature_extractor)
    total_gen_loss = gan_loss + (100 * l1_loss) + (10 * perceptual)
    return total_gen_loss, gan_loss, l1_loss

def discriminator_loss(disc_real_output, disc_generated_output):
    real_loss = loss_object(tf.ones_like(disc_real_output), disc_real_output)
    generated_loss = loss_object(tf.zeros_like(disc_generated_output), disc_generated_output)
    total_disc_loss = real_loss + generated_loss
    return total_disc_loss

generator_optimizer = tf.keras.optimizers.Adam(2e-4, beta_1=0.5)
discriminator_optimizer = tf.keras.optimizers.Adam(2e-4, beta_1=0.5)

```

Figure 5.1. Training Setup.

The model trains on 50 epochs. During each training iteration the generator generates a new image from the sketch, and then the discriminator evaluates the pair of real and generated images. The gradients for each of the networks are calculated and used separately by their own optimizers, thus sustaining their adversarial game that pushes them to progress. At each batch count divisible by 10, the loss values for the current generator and discriminator are outputted to the console, while the average losses on the train and validation datasets are reported at the end of each epoch.

The total generator loss is a combination of three losses, each having their respective

weights to balance the influence over the entire objective. The adversarial loss gets a 1 weight, ensuring the generator is still in an adversarial relationship with the discriminator. The L1 pixel loss has the weight of 100, thus being the main structural factor — it is the loss which ensures that the generator is generating something similar to the target dataset. Perceptual loss against the deep features of VGG19 gets a weight of 10, balancing the structural guidance from L1.

5.2 Validation Process and Testing Process

After every training epoch, the model is assessed using a held-out validation set consisting of 20 percent of the entire data set. In this case, the generator is run in inference mode, where dropout and batch normalization work as they would under the test conditions, and the average generator loss is calculated using all validation batches. This provides a good way of knowing whether the model has made any improvement or has just overfit the training data.

Qualitative tests involve generation of images in pairs of the input sketch, target image, and output image for 10 randomly chosen validation samples. All pixel values in the output image are adjusted to their normalized form in the range of [0, 1]. These kinds of visualization provide good insight into how the generator can handle the texture and structure of real images.

```
@tf.function
def train_step(input_image, target):
    with tf.GradientTape() as gen_tape, tf.GradientTape() as disc_tape:
        gen_output = generator(input_image, training=True)

        disc_real_output = discriminator([input_image, target], training=True)
        disc_generated_output = discriminator([input_image, gen_output], training=True)

        gen_loss, gan_loss, l1_loss = generator_loss(disc_generated_output, gen_output, target)
        disc_loss = discriminator_loss(disc_real_output, disc_generated_output)

        generator_gradients = gen_tape.gradient(gen_loss, generator.trainable_variables)
        discriminator_gradients = disc_tape.gradient(disc_loss, discriminator.trainable_variables)

        generator_optimizer.apply_gradients(zip(generator_gradients, generator.trainable_variables))
        discriminator_optimizer.apply_gradients(zip(discriminator_gradients, discriminator.trainable_variables))

    return gen_loss, disc_loss
```

Figure 5.2. Validation & Testing Setup.

```

@tf.function
def val_step(input_image, target):
    gen_output = generator(input_image, training=False)
    gen_loss, _ = generator_loss(discriminator([input_image, gen_output], training=False), gen_output, target)
    return gen_loss

def train(dataset, val_dataset, epochs):
    for epoch in range(epochs):
        print(f"\nEpoch {epoch + 1}/{epochs}")
        gen_losses = []
        disc_losses = []

        for n, (input_image, target) in enumerate(dataset):
            gen_loss, disc_loss = train_step(input_image, target)
            gen_losses.append(gen_loss)
            disc_losses.append(disc_loss)
            if n % 10 == 0:
                print(f"Batch {n}: Gen Loss = {gen_loss.numpy()}, Disc Loss = {disc_loss.numpy()}")

        print(f"Epoch {epoch + 1} Training: Avg Gen Loss = {np.mean(gen_losses)}, Avg Disc Loss = {np.mean(disc_losses)}")

        val_losses = []
        for input_image, target in val_dataset:
            val_loss = val_step(input_image, target)
            val_losses.append(val_loss)
        print(f"Epoch {epoch + 1} Validation: Avg Gen Loss = {np.mean(val_losses)}")

train(train_dataset, val_dataset, epochs=50)

```

Figure 5.3. Validation & Testing Setup

5.3 Results

Model convergence was uniform throughout the 50 epochs with the reduction of the generator loss and discriminator loss, as shown below:

Epoch 1: Generator Loss = 61.20; Discriminator Loss = 1.71; Validation Loss = 49.92, clearly showing an indication of the beginning of adversarial inconsistency.

Epoch 3: The generator loss reduces to 38.45, meaning a decrease of about 37 percent within just two epochs because of the L1 loss being multiplied by 100, leading the generator towards structure resemblance.

Epoch 50: During this epoch, the generator loss reduced to 13.36, discriminator loss converged at 1.08, while the validation loss was also reduced to 32.06.

Convergence & Generalization: The consistency of the discriminator loss during training epochs was in the range of 1.0 to 1.2, showing good balance between the two processes and also demonstrating generalization capabilities through train-validation loss reduction during the epochs.

Qualitative observations: Based on qualitative observation of the generated images using 10 random validation images, it has been observed that the structure and

symmetry of the face were maintained quite efficiently by U-Net skip connections. Colors including skin, hair and the background were reproduced very effectively and it can be said that the network efficiently learned low frequency mapping of the image. Features of VGG19 Block5_Conv2 helped improve the semantics of details such as eyes, nose, and hairline, which cannot have been achieved with just L1 loss.

```
Epoch 1/50
Batch 0: Gen Loss = 73.58887481689453, Disc Loss = 1.9172000885009766
Epoch 1 Training: Avg Gen Loss = 61.202064514160156, Avg Disc Loss = 1.714606523513794
Epoch 1 Validation: Avg Gen Loss = 49.9174690246582

Epoch 2/50
Batch 0: Gen Loss = 50.12726593017578, Disc Loss = 1.6017847061157227
Epoch 2 Training: Avg Gen Loss = 45.555267333984375, Avg Disc Loss = 1.1459600925445557
Epoch 2 Validation: Avg Gen Loss = 39.05696487426758

Epoch 3/50
Batch 0: Gen Loss = 39.20075988769531, Disc Loss = 1.1631526947021484
Epoch 3 Training: Avg Gen Loss = 38.450199127197266, Avg Disc Loss = 1.1828628778457642
Epoch 3 Validation: Avg Gen Loss = 34.21597671508789

Epoch 4/50
Batch 0: Gen Loss = 35.88298416137695, Disc Loss = 1.317540168762207
Epoch 4 Training: Avg Gen Loss = 35.75533676147461, Avg Disc Loss = 1.187007188796997
Epoch 4 Validation: Avg Gen Loss = 32.53214645385742

Epoch 5/50
Batch 0: Gen Loss = 37.54705810546875, Disc Loss = 0.9763821363449097
...
Epoch 50/50
Batch 0: Gen Loss = 15.02348518371582, Disc Loss = 0.8043128252029419
Epoch 50 Training: Avg Gen Loss = 13.361335754394531, Avg Disc Loss = 1.0761632919311523
Epoch 50 Validation: Avg Gen Loss = 32.06052780151367
Output is truncated. View as a scrollable element or open in a text editor. Adjust cell output settings...
```

Figure 5.4. Training The Model



Figure 5.5. Final Output

CHAPTER 6

CONCLUSION, LIMITATIONS AND FUTURE SCOPE

6.1 Conclusion

In this regard, the model aimed at exploring the possibilities of deep learning in terms of filling the gap between a rough sketch and photorealistic image of a person's face. Based on the results, it should be noted that there is a positive answer to the question posed above as this model was able to accomplish this goal. It should be noted that the use of the architecture of the pix2pix model with the perceptual loss function based on VGG19 made it possible to train this model to transfer sketches into photorealistic images in an appropriate way without matching pixel-by-pixel operation.

In particular, it was found that the U-Net architecture used in the generator was adequate enough for the task at hand. In turn, the use of the cross-layer skip connection was justified due to the opportunity to pass the necessary spatial information which gets lost during the propagation through the deep bottleneck. Thus, the generator became able to create semantically correct facial attributes. Moreover, instance normalization made it possible to stabilize the activation values in different areas of input faces. Finally, PatchGAN discriminator helped to learn the proper composition and details.

In conclusion, what this paper proves is that through using the right balance of adversarial, structure, and perceptual objectives in paired image-to-image translation, one can create something that looks impressive from a technical perspective while simultaneously being practically relevant. The method used here is far from flawless since one notices that there is still a blurring of high frequencies that occurs at times, but nevertheless this is a solid base to build upon in the future.

6.2 Challenges and Limitations

One of the most persistent challenges in training any GAN is keeping the generator and discriminator in a productive equilibrium. If the discriminator learns too quickly, the generator receives gradients that are too uninformative to improve from. If the generator pulls ahead, the discriminator loses its ability to provide a useful learning signal. In this model, the choice of Adam with $\beta_1 = 0.5$ and a learning rate of $2e-4$ helped manage this tension, but the early epochs — particularly epoch 1, where discriminator loss opened at 1.71 — showed clear signs of instability before the two networks settled into a more balanced dynamic.

This is because the model learns the one-to-one mapping relationship between the two datasets of sketches and photographs. As a result, the model usually becomes sensitive to the variations in terms of the style, lines, and diversity of the faces used in the drawing. When the sketches exceed the distribution learned by the model from the training data, such as stylized sketches, cartoons, or even sketches having a different proportion, the outcomes will definitely be suboptimal. The fact that the output sketches exhibit almost the same uniformly-colored blue-gray background points to the possible overfitting issue to the background in the training data. Training a pix2pix model with a VGG19 perceptual loss is computationally demanding. The inclusion of the pre-trained VGG19 network as a feature extractor added meaningful overhead to each training step, since forward passes through the VGG19 network must be computed for both the target and generated images at every batch update. With a batch size of 16 and 50 epochs, this added up to a substantial training time, and would scale poorly with larger datasets or higher input resolutions without access to multi-GPU infrastructure.

The difference between training generator loss (13.36) and validation generator loss (32.06) at epoch 50, although nothing to be overly concerned about, does imply that there is some overfitting towards the training data. It is natural to have this overfitting in cases where there is an effort made in creating pairs of corresponding images; however, the result still implies the necessity of applying some regularizing techniques to ensure that the validation performance approaches the training one.

6.3 Future Scope

Although the current model was created for the specific task of generating faces from sketches, its architecture is applicable to various problems involving two images, input and output images. It means that we can use it, for example, to create architectural drawings based on sketches, create medical anatomy based on sketches, photorealistic product pictures based on its sketch image, or even generate satellite pictures based on maps. They would all need different datasets for training, but they can use the same model and loss function.

The simplest thing that we can do to improve this result is increasing output resolution to either 512x512 or 1024x1024. In order to do so, we need better equipment or, alternatively, make some architectural changes in our network. This change will lead to another approach to solving the issue of low-resolution and loss of information that is inevitable in the case of the current model.

Finally, future research may use more appropriate metrics to evaluate sketch-to-photo generators' results other than losses. Metrics, such as Fréchet Inception Distance (FID), Structural Similarity Index (SSIM), and Peak Signal to Noise Ratio (PSNR) can serve as an objective way to compare different models' results.

Moreover, performing user evaluations wherein humans judge how realistic and personality-sketched is the created image will give insights that are not provided via any other numerical measure.

In the transition process to deployment of the proposed model, it could be interesting to explore methods of neural network compression such as knowledge distillation, quantization, and pruning. Through this method, we can deploy the proposed model in an efficient way to provide real-time translation, making it interactive and providing photorealistic results to those who use it for painting.

CHAPTER 7

SOCIAL IMPACT

However, the potential of transforming such a sketch into an image may prove highly practical. While it is evident that the particular technological advancement would certainly play a major role in the future, it is also necessary to admit that its significance may depend greatly on what particular use it would receive.

In connection with one of the possible ways of implementing this kind of technology, it would be appropriate to point to the use of this technology in the field of forensics. In this regard, the use of such forensic sketches had long proved their importance in establishing criminal persons by the police. However, the problem was always the quality of the artist or of the witnesses' accounts. It would thus be advantageous for the police to use the suggested system in establishing the identity of the suspects or lost persons.

It is just the reason why the innovative technology may prove to be quite dangerous. Thus, the use of such a drawing technology would entail a certain risk, since:

- Privacy breach: The use of this algorithm in order to generate a depiction of the face of a living individual based solely on a verbal description of the face or an individual's personal recollections regarding the appearance of a certain face is a breach of the individual's privacy.
- The chances of mistaken identity: In case this technology is utilized in forensic research, the generated facial depiction could be misused in the research process by making wrong decisions and even accusations about the innocence of an individual based on bias present in the training data set.
- The potential for face generation to be exploited: This kind of technology is already at the stage where it can easily generate content. Consequently, it would be relatively easy to exploit that technology in order to generate depictions of real people's faces.

It should be considered ethically important to implement any technology in such a way that allows it to bring some social benefit. In the case at hand, ethical considerations could include transparent functioning of the system, probabilistic nature of the work, need for human decision-making in cases such as forensics, and regular checking of the data set.

REFERENCES

- [1] 10.53070-bbd.1429596-3698619.pdf: S. Altun Güven, E. Şahin, and M. F. Talu, "Image-to-Image Translation with CNN Based Perceptual Similarity Metrics," *Journal of Computer Science (Anatolian Science)*, vol. 9, no. 1, pp. 84-98, 2024.
- [2] 1406.2661v1.pdf: I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative Adversarial Nets," *Advances in Neural Information Processing Systems (NIPS)*, 2014.
- [3] 1505.04597v1.pdf: O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional Networks for Biomedical Image Segmentation," in *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*, pp. 234–241, 2015.
- [4] 1603.08155v1.pdf: J. Johnson, A. Alahi, and L. Fei-Fei, "Perceptual Losses for Real-Time Style Transfer and Super-Resolution," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016.
- [5] 1611.07004v3.pdf: P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-Image Translation with Conditional Adversarial Networks" (also known as Pix2Pix), in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [6] 1703.10593v7.pdf: J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks" (also known as CycleGAN), in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017.
- [7] 1711.11585v2.pdf: T.-C. Wang, M.-Y. Liu, J.-Y. Zhu, A. Tao, J. Kautz, and B. Catanzaro, "High-Resolution Image Synthesis and Semantic Manipulation with Conditional GANs" (also known as Pix2PixHD), in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [8] 1905.01270v2.pdf: H.-Y. Lee, H.-Y. Tseng, Q. Mao, J.-B. Huang, Y.-D. Lu, M. Singh, and M.-H. Yang, "DRIT++: Diverse Image-to-Image Translation via Disentangled Representations," *International Journal of Computer Vision*

(IJCV), 2019.

[9] 1909.08313v3.pdf: R. Liu, Q. Yu, and S. Yu, "An Unpaired Sketch-to-Photo Translation Model," *arXiv preprint arXiv:1909.08313*, 2019.

[10] 2012.09290v2.pdf: B. Liu, K. Song, and A. Elgammal, "Sketch-to-Art: Synthesizing Stylized Art Images From Sketches," *arXiv preprint arXiv:2002.12888*, 2020.

[11] 2211.09800v2.pdf: T. Brooks, A. Holynski, and A. A. Efros, "InstructPix2Pix: Learning to Follow Image Editing Instructions," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.

[12] 2402.17624v1.pdf: N. Tumanyan, M. Geyer, S. Bagon, and T. Dekel, "Plug-and-Play Diffusion Features for Text-Driven Image-to-Image Translation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.

[13] 2404.19265v2.pdf: Z. Li, B. Guan, Y. Wei, and Y. Zhou, "Mapping New Realities: Ground Truth Image Creation with Pix2Pix Image-to-Image Translation," 2024.

[14] 2408.06000v1.pdf: (The provided source for this file [405–424] contains references only; however, the content relates to Self-Supervised Consistency Regularization and Style Transfer for I2I translation).

[15] 2408.14270v1.pdf: L. Zhou and G. Li, "MITIA: Misalignment-Invariant Training for Unpaired Medical Image-to-Image Translation," *Medical Physics*, 2024.

[16] 2502.06895v1.pdf: W. Jiangtao, N. I. R. Ruhaiyem, and F. Panpan, "A Comprehensive Review of U-Net and Its Variants: Advances and Applications in Medical Image Segmentation," 2025.

[17] 2505.16310v1.pdf: (This document covers generative models for Cityscapes and Pap-Smear image generation using U-Net and GAN architectures [517–523]).

[18] 2601.04785v1.pdf: X. Qiu, Y. Dai, X. Tan, S. Li, F. Sun, L. Gan, and L. Liu,

"SRU-Pix2Pix: A Squeeze-and-Excitation Residual U-Net for Medical Image-to-Image Translation," 2026 [524–526].

[19] 2602.16664v1.pdf: J. Liu, F. Petersen, Y. Gao, Y. Zhang, H. Kim, A. S. Chaudhari, Y. Sun, S. Ermon, and S. Gatidis, "Unpaired Image-to-Image Translation via a Self-Supervised Semantic Bridge" (SSB), 2026.

[20] 2604.12805v1.pdf: F. Tan, H. Yang, Q. Duan, Q. Xie, D. Meng, and K. Ye, "TLEQ-CycleGAN: Rotation Equivariant Image-to-Image Translation," 2026.

[21] 3606694.pdf: D. Donoso and J. M. Saavedra, "Survey on Sketch-to-photo Translation," *ACM Computing Surveys*, vol. 56, no. 1, art. 22, pp. 1–25, 2023.

[22] Image_to_Image_Translation_Based_on_Differential_I.pdf: X. Zhao, H. Yu, and H. Bian, "Image to Image Translation Based on Differential Image Pix2Pix Model," *Computers, Materials & Continua*, vol. 77, no. 1, pp. 189–198, 2023.

5%

Overall Similarity

The combined total of all matches, including overlapping sources, for each database.

Filtered from the Report

- Bibliography
- Quoted Text
- Cited Text
- Small Matches (less than 8 words)

None

Submission ID: trn:old-9832-140478485

Match Groups

- 48 Not Cited or Quoted: 5%
Matches with neither in-text citation nor quotation marks
- 0 Missing Quotations: 0%
Matches that are still very similar to source material
- 0 Missing Citation: 0%
Matches that have quotation marks, but no in-text citation
- 0 Cited and Quoted: 0%
Matches with in-text citation present, but no quotation marks

Top Sources

- 3% Internet sources
- 2% Publications
- 4% Submitted works (Student Papers)

Integrity Flags

0 Integrity Flags for Review

No suspicious text manipulations found.

Our system's algorithms look deeply at a document for any inconsistencies that would set it apart from a normal submission. If we notice something strange, we flag it for you to review.

A Flag is not necessarily an indicator of a problem. However, we do recommend you focus your attention there for further review.

Submission ID: trn:old-9832-140478485

*0% detected as AI

AI detection includes the possibility of false positives. Although some text in this submission is likely AI generated, scores below the 20% threshold are not surfaced because they have a higher likelihood of false positives.



Caution: Review required.

It is essential to understand the limitations of AI detection before making decisions about a student's work. We encourage you to learn more about Turnitin's AI detection capabilities before using the tool.

Disclaimer

Our AI writing assessment is designed to help educators identify text that might be prepared by a generative AI tool. Our AI writing assessment may not always be accurate (i.e., our AI models may produce either false positive results or false negative results), so it should not be used as the sole basis for adverse actions against a student. It takes further scrutiny and human judgment in conjunction with an organization's application of its specific academic policies to determine whether any academic misconduct has occurred.