

CS424 – Group Project Report

Akeela Darryl Fattha
akeelaf.2022@scis.smu.edu.sg

Fadhel Erlangga Wibawanto
fadhelew.2022@scis.smu.edu.sg

Tan Zhi Rong
zhirong.tan.2022@scis.smu.edu.sg

Grace Angel Bisawan
gbisawan.2022@scis.smu.edu.sg

Lee Jia Heng
jiaheng.lee.2023@scis.smu.edu.sg

Abstract

This report provides an overview of Cycle-Consistent Adversarial Networks (CycleGAN), a technique for unpaired image-to-image translation. We explore the architecture, key innovations, applications, and limitations of CycleGAN models in computer vision tasks.

Contents

1	Task 1	3
1.1	Introduction	3
1.1.1	Configuration & Hyperparameters	3
1.2	Architecture	3
1.2.1	Generator	3
1.2.2	Discriminator	4
1.3	Data Preparation	4
1.3.1	Augmentations	4
1.3.2	Data Split	5
1.4	Methodology	5
1.4.1	Staged Training	5
1.4.2	Image Pooling	6
1.5	Loss Functions	6
1.5.1	Discriminator Losses	6
1.5.2	Generator Losses	7
1.6	Conclusion	7
1.7	Appendix	7
2	Task 2	11
2.1	Introduction	11
2.1.1	Configuration	11
2.2	Architecture	11
2.2.1	Generator	11
2.2.2	Discriminator	11
2.3	Data Preparation	11
2.3.1	Collection	11
2.3.2	Pre-Processing	13
2.3.3	Data Split	13
2.4	Loss Functions	13
2.4.1	Discriminator Losses	13
2.4.2	Generator Losses	14
2.5	Methodology	15
2.6	Metrics	15
2.7	Conclusion	15
2.8	Appendix	15
2.8.1	Randomly Paired vs Supervised Paired Dataset Preference	15

1. Task 1

1.1 Introduction

We attempt to create a CyleGAN using novel and modern techniques to improve the performance of the model to convert cartoon faces to realistic faces, and vice-versa. We will be using improved blocks, various attention mechanisms and normalisation techniques in our generators, as well as additional loss functions to better guide the models to learn the mapping between the two domains. We opted not to use any augmentations to our images, and we use an image size of 256x256 as our input/output image size.

With our enhancements, we achieved the following results:

	FID	IS	GMS	Avg GMS
Raw to Cartoon	46.70193	2.41466 ± 0.24392	4.39783	3.85666
Cartoon to Raw	39.83461	3.62382 ± 0.23992	3.31548	

Further supplementary details can be found in the Appendix 1.7.

1.1.1 Configuration & Hyperparameters

All optimisers used Adam with a Cosine Annealing Warm Restarts scheduler. The following hyperparameters and loss weights were used for our best results:

- Random seed: 42
- $Scheduler_{t-mult}$: 2
- λ_{gan} : 1.0
- Learning rate: 2e-4
- $Scheduler_{min-lr}$: 1e-5
- λ_{gp} : 10.0
- Betas: (0.5, 0.999)
- λ_{cyc} : 10.0
- λ_{id} : 5.0
- λ_{fm} : 5.0
- $Scheduler_{t-0}$: 10

1.2 Architecture

1.2.1 Generator

Our generator incorporates several modern attention mechanisms and architectural techniques including Convolutional Block Attention Module (CBAM), Squeeze-and-Excitation (SE) blocks, Self-Attention mechanisms, Spatial Attention Blocks and Channel Attention Blocks for improved feature representation, generating higher-quality images with better detail preservation, style consistency, and structural coherence. Some of our key design choices are listed below:

- **Instance Normalization:** Used throughout the network instead of batch normalization, as it has been shown to produce better results for style transfer and image-to-image translation by normalizing each instance independently.
- **Reflection Padding:** Applied before convolutions to reduce boundary artifacts that can appear in generated images, particularly important for maintaining realistic edges.
- **Mixed + Improved Residual Blocks:** The strategic placement of different residual block types allows the network to benefit from complementary approaches to feature transformation. They also exhibit the following characteristics and sub-blocks:
 - **Channel Attention:** Models interdependencies between channels using both max and average pooling operations.
 - **Spatial Attention:** Focuses on important spatial regions by creating attention maps from channel-wise statistics.
 - **Skip Connection:** Allows gradients to flow more easily through the network, improving convergence and enabling the model to learn more complex mappings.

– All blocks benefit from SE attention for channel recalibration

- **Squeeze-and-Excitation Attention:** Used in Mixed Residual Blocks, it recalibrates channel-wise feature responses by explicitly modeling interdependencies between channels, allowing the network to selectively emphasize informative features.
- **Enhanced Self-Attention:** Positioned after the residual blocks, it helps ensure global coherence in the generated images by modeling long-range dependencies that convolutional operations cannot capture efficiently. This is particularly valuable for maintaining structural integrity in facial regions.

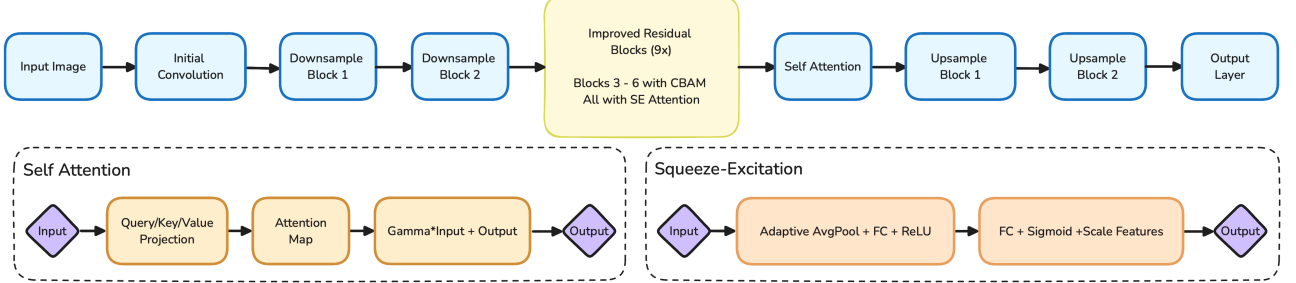
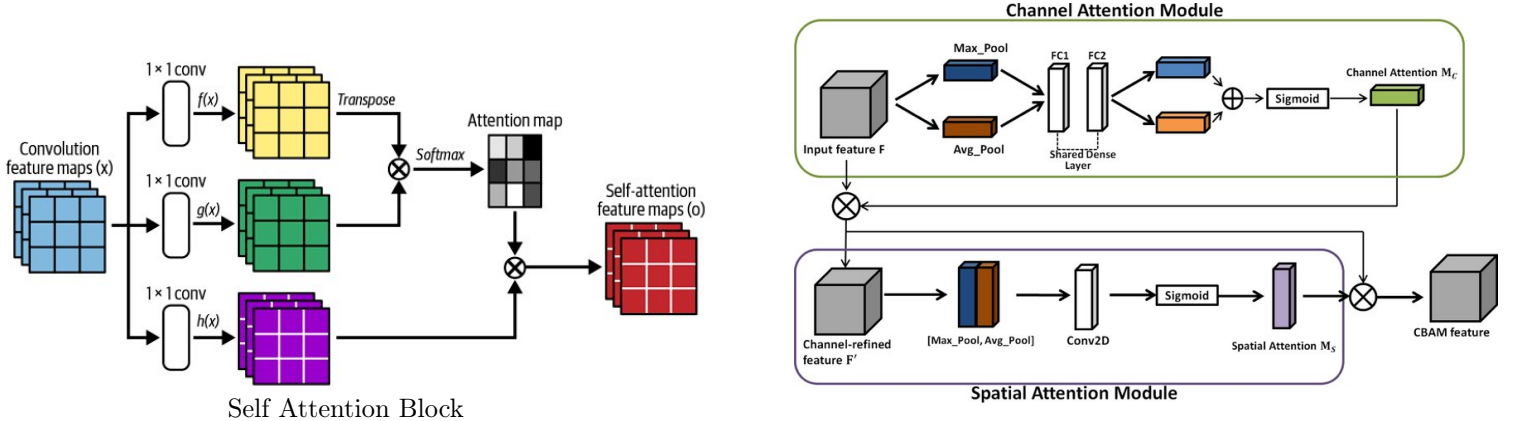


Figure 1.1: Generator Architecture



1.2.2 Discriminator

Our discriminator uses a basic sequential architecture with convolutional layers that doubles in number of channels 3 times, from 64 to 512 and downsizing the image size by half 4 times to 16x16. Each convolution output is passed to a spectral norm layer, instance normalisation and a leaky relu activation function. We then branched into a Avg Pool into a FC layer for global classification and a PatchGAN block for local classification.

Given that we have used a 256x256 image, the feature map will have size of $[B, 512, 16, 16]$, the patchGAN will have size of $[B, 1, 16, 16]$. A larger patch GAN and feature map allows us to capture more details of specific features in certain areas which allows the discriminator to provide more useful information to the generator.

1.3 Data Preparation

1.3.1 Augmentations

We started with resizing images to 128x128, with a batch size of 16 before we hit a memory limit. We also applied no augmentations to the images, as a baseline to compare against later on when we add augmentations.

After adding various minor augmentations like affines, color jitters, flips, light gaussian blurs and posterising, we found that any combination of them resulted in worse metrics overall. From visual inspection, we found that the generated raw images were more noisy than without augmentations, which could contribute to the worse

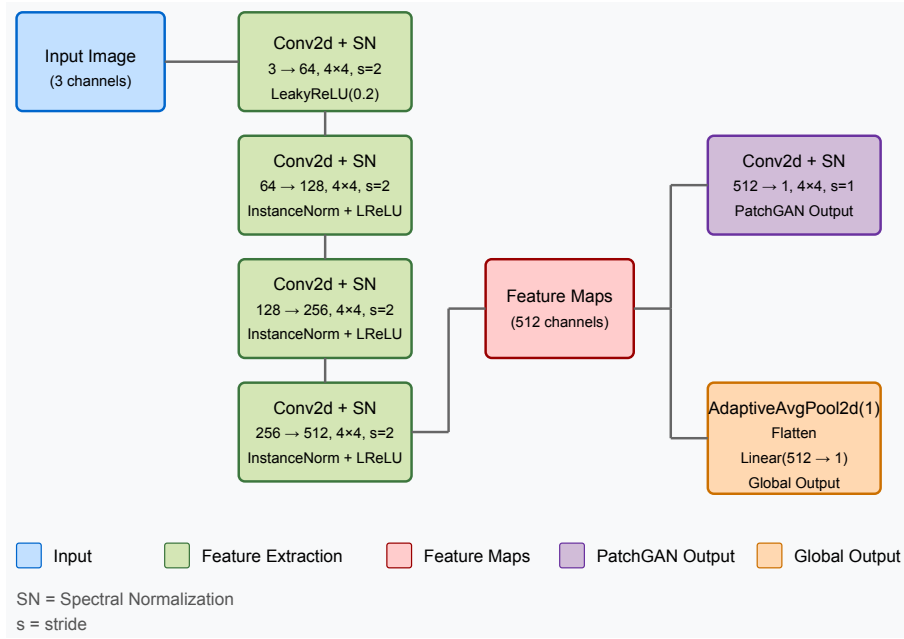


Figure 1.2: Task 1: Discriminator Architecture

metrics. Furthermore, the generated cartoon images were a lot more flat than the target images, losing a lot more of the finer details and textures. Lastly, affines and any crops were clearly visible in the generated images when they otherwise should not be.

Thus, in both cases, we believe that adding the augmentations were encouraging the model to over-process the style conversions, which likely caused the reduction in metrics. Finally our best result was achieved with no augmentations, and upscaling the images to 256x256 with a batch size of 3. Some sample results are shown in the below Figures.

1.3.2 Data Split

We worked with simple 99/1 split on training/validation data since we found that we needed as much data as possible for training the model to achieve our best results.

1.4 Methodology

1.4.1 Staged Training

Our best results came from a multi-staged supervised training approach. The idea is inspired from an application on Sketch-to-Manga colorization by Zhang et al. (2018)¹, where they utilised 2 generators using the same architecture, intended for them to have different focuses. Our methodology follows a "greedy" approach seen in **Algorithm 1** where we iteratively save checkpoints of our models, measure their metrics, and we keep the best performing weights. We then use these weights as a starting point for the next stage of training.

Our novelty stems from how we decrease the batch size, number of epochs, and evaluation threshold proportionally in each iterative stage, while using newly randomised seeds and maintaining learning rate so that our ratio of batch size to learning rate decreases². This allows us to gradually refine the model's performance and achieve better results as each iteration encourages the model to learn finer details. Our intention with newly randomised seeds was to allow the optimisers to start on a different location in the loss optimisation graph, which may force it to hit a new local minima.

¹Zhang, L., Li, C., Wong, T.-T., Ji, Y., & Liu, C. (2018). Two-stage sketch colorization. ACM Transactions on Graphics, 261–261. <https://ttwong12.github.io/papers/colorize/colorize.pdf>

²He, F., Liu, T., & Tao, D. (2019). Control Batch Size and Learning Rate to Generalize Well: Theoretical and Empirical Evidence. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d' Alché-Buc, E. Fox, & R. Garnett (Eds.), Advances in Neural Information Processing Systems (Vol. 32). Retrieved from https://proceedings.neurips.cc/paper_files/paper/2019/file/dc6a70712a252123c40d2adba6a11d84-Paper.pdf

Stage	Batch Size	Epochs	Eval Threshold	Torch Seed	D_A	D_B	G_{AB}	G_{BA}
Stage 1	3	100	10	1722426635407800	–	–	–	–
Stage 2	2	50	10	1812572933112200	100	100	100	100
Stage 3	1	30	3	–	130	150	150	130

Table 1.1: Configuration of hyperparameters for each stage. Model weights are based on epoch when checkpoint was created.

Algorithm 1 Basic Training Loop Single Stage

```

procedure TRAINMODEL( $D_A^w, D_B^w, G_{AB}^w, G_{BA}^w, batchSize, epochs, evalThreshold$ ) ▷ weights for models
   $D_A \leftarrow D_A^w$ 
   $D_B \leftarrow D_B^w$ 
   $G_{AB} \leftarrow G_{AB}^w$ 
   $G_{BA} \leftarrow G_{BA}^w$ 
   $n \leftarrow 1$ 
   $W\_history \leftarrow [ ]$ 
  while  $n \neq epochs$  do ▷ Run basic training loop
     $Train \ D_A, D_B, G_{AB}, G_{BA}$  with  $Domain_A, Domain_B, batchSize$ 
    if  $n \bmod evalThreshold = 0$  then
       $W\_history.insert(D_A^w, D_B^w, G_{AB}^w, G_{BA}^w)$ 
    end if
     $n \leftarrow n + 1$ 
  end while
   $bestGMS \leftarrow inf$ 
   $bestW \leftarrow D_A^w, D_B^w, G_{AB}^w, G_{BA}^w$ 
  for  $w_e \in W$  do ▷ Evaluate each weight with test sets
     $D_A \leftarrow D_A^{w_e}$ 
     $D_B \leftarrow D_B^{w_e}$ 
     $G_{AB} \leftarrow G_{AB}^{w_e}$ 
     $G_{BA} \leftarrow G_{BA}^{w_e}$ 
     $epochGMS \leftarrow Compute \ GMS$ 
    if  $epochGMS < bestGMS$  then
       $bestGMS \leftarrow epochGMS$ 
       $bestW \leftarrow w_e$ 
    end if
  end for
  return  $bestW$  ▷ Returns the best weights
end procedure

```

1.4.2 Image Pooling

An "Image Pool" was also used in training with a buffer of 50 images. This pool functions as a buffer to store previously generated images, which helps to reduce the risk of mode collapse where the Discriminator becomes too powerful, and prevents the Generator from learning anything new. The pool is updated with new generated fake images during training, and once the threshold is reached, we randomly sample from the available pool instead of using the latest generated image.

1.5 Loss Functions

1.5.1 Discriminator Losses

We used a combination of loss functions beyond the standard adversarial loss to improve the performance of our discriminator. Our overall loss function for the discriminator is as follows:

$$L_D = L_{GAN} + \lambda_{gp} L_{gp} \quad (1.1)$$

Gradient Penalty – L_{gp}

We also use a gradient penalty, inspired from Wasserstein GAN (WGAN)³ in order to prevent the discriminator from learning too fast. Doing so avoids problems like mode collapse with vanishing gradients, and ensures that the generator is always learning, even if it is not performing well enough to "beat" the discriminator yet, thus improving stability of training.

$$\text{Gradient Penalty Loss} = 0.7 \cdot \mathbb{E} \left[(\|\nabla_{\hat{x}} D_{\text{patch}}(\hat{x})\|_2 - 1)^2 \right] + 0.3 \cdot \mathbb{E} \left[(\|\nabla_{\hat{x}} D_{\text{global}}(\hat{x})\|_2 - 1)^2 \right]$$

Relativistic – L_{GAN}

We extend the typical adversarial loss by using a relativistic loss, which compares the logits of discriminating the real and fake images in a relative manner on a global scope with Binary Cross Entropy Loss. Furthermore, we took inspiration from PatchGAN⁴ which whether the images is real or fake based on smaller patches, a local scope. Thus, we have another source of information for determining real or fake, which is computed with MSE loss. The equation for our relativistic loss for our "real" side is as follows, and the opposite will be applied for the "fake" side i.e. $D_{\text{real}} <=> D_{\text{fake}}, \text{True} <=> \text{False}$:

$$L_{GAN} = 0.7 \cdot \text{MSE}(D_{\text{real}} - \bar{D}_{\text{fake}, \text{True}})_{\text{patch}} + 0.3 \cdot \text{BCE}(D_{\text{real}} - \bar{D}_{\text{fake}, \text{True}})_{\text{global}} \quad (1.2)$$

1.5.2 Generator Losses

The generator also makes use of the **Relativistic** loss, which is an extension of the standard adversarial GAN loss. Our overall loss function for the generator is as follows:

$$L_G = L_{GAN} + \lambda_{\text{cyc}} L_{\text{cyc}} + \lambda_{\text{id}} L_{\text{id}} + \lambda_{\text{feat}} L_{\text{feat}} \quad (1.3)$$

Cycle Consistency – L_{cyc}

As we would like the fake images to be able revertible to its original image i.e. the pixels are identical, we maintain using L1 Loss for our Cycle Consistency.

Identity – L_{id}

Similarly, we also use L1 Loss for our Identity loss since images in the same domain should not be altered.

Feature – L_{feat}

To encourage perceptual similarity between generated Pokémon and their corresponding animal inputs, we incorporate a Feature Loss (L_{feat}) based on intermediate activations from our discriminator. While perceptual loss is often computed using pre-trained networks like VGG, these models are trained on datasets such as ImageNet, which do not share the stylistic or semantic characteristics of our target domain. Instead, we extract features from our task-specific discriminator, which learns domain-relevant representations during adversarial training. We use the L1 norm to compute the distance between features of the generated and target images to focus on accurate reconstruction, as it is more robust to outliers and tends to preserve texture and structure better than L2, which can overly penalize large deviations and lead to smoother, less detailed outputs.

1.6 Conclusion

In summary, our work extends the typical CycleGAN framework with significant enhancements to both architecture and training methodology. By integrating advanced attention mechanisms and refined normalization techniques (see 1.2.1), and adopting a multi-staged training approach with strategically adjusted hyperparameters (see 1.4.1), we achieved significantly improved image quality, from METRIC1 down to METRIC2. The combined use of gradient, adversarial, cycle consistency, identity, and feature losses (see 1.5) further stabilised training and perceptual realism.

1.7 Appendix

³Arjovsky, M., Chintala, S., & Bottou, L. (2017, July 17). Wasserstein Generative Adversarial networks. PMLR. <https://proceedings.mlr.press/v70/arjovsky17a.html>

⁴Isola, P., Zhu, J., Zhou, T., & Efros, A. A. (2016). Image-to-Image Translation with Conditional Adversarial Networks. arXiv (Cornell University). <https://doi.org/10.48550/arxiv.1611.07004>

Stage / Direction	FID ↓	IS ↑	GMS ↓	Avg GMS ↓
Stage 1				
Raw to Cartoon	46.70193	2.41466 ± 0.24392	4.39783	3.85666
Cartoon to Raw	39.83461	3.62382 ± 0.23992	3.31548	
Stage 2				
Raw to Cartoon	31.6551	2.2699 ± 0.20	3.7344	3.42423
Cartoon to Raw	35.1988	3.6300 ± 0.22	3.1139	
Stage 3				
Raw to Cartoon	TBC	TBC ± 0.18	TBC	TBC
Cartoon to Raw	TBC	TBC ± 0.21	TBC	

Table 1.2: Evaluation metrics across stages, best results taken from each stage



Figure 1.3: Original (resized 256x256). Top row: Raw, Bottom row: Cartoon



Figure 1.4: 128x128 Avg GMS: 4.47829



Figure 1.5: 256x256, Stage 1: Batch Size 3, Avg GMS: 3.72523



Figure 1.6: 256x256, Stage 2: Batch Size 2, Avg GMS: 3.42415

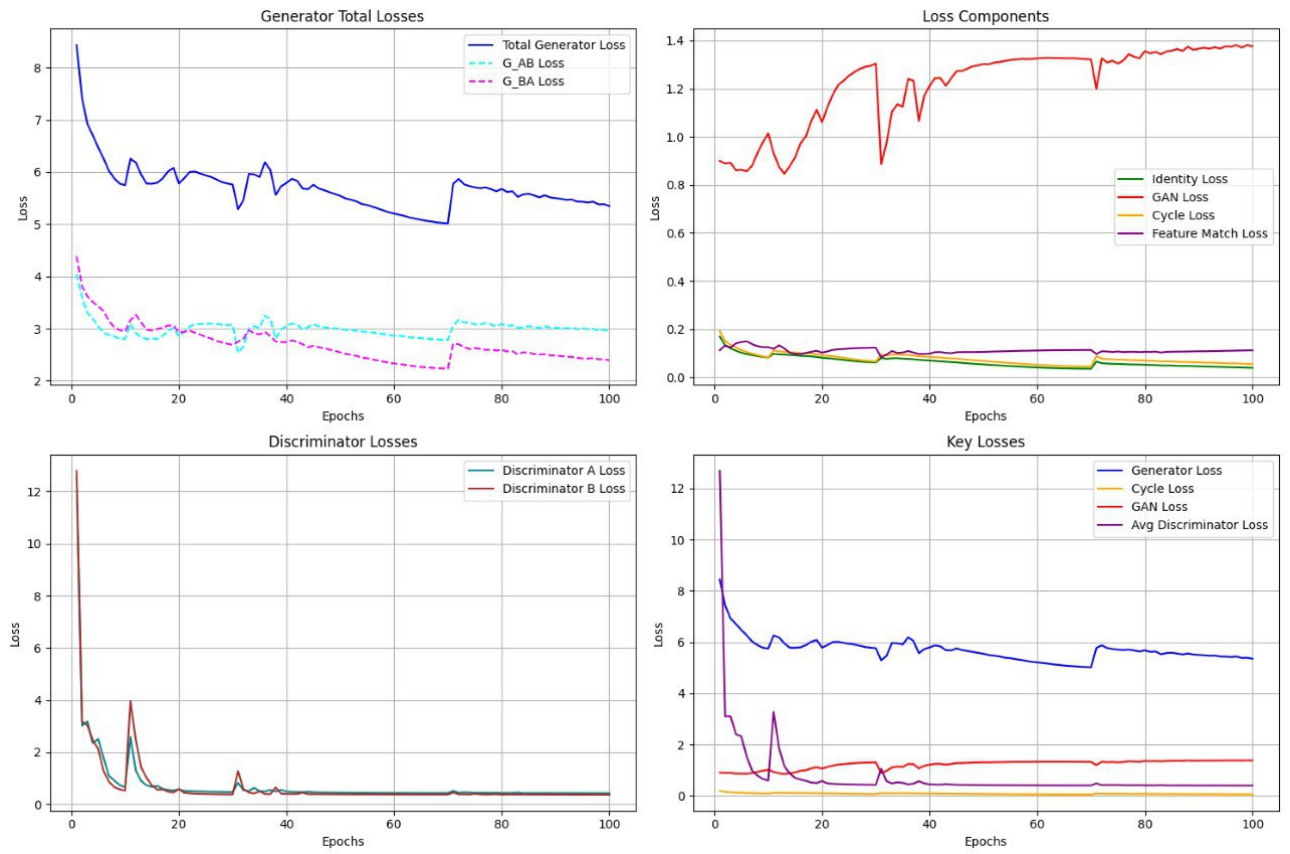


Figure 1.7: Loss curve for Stage 1 training

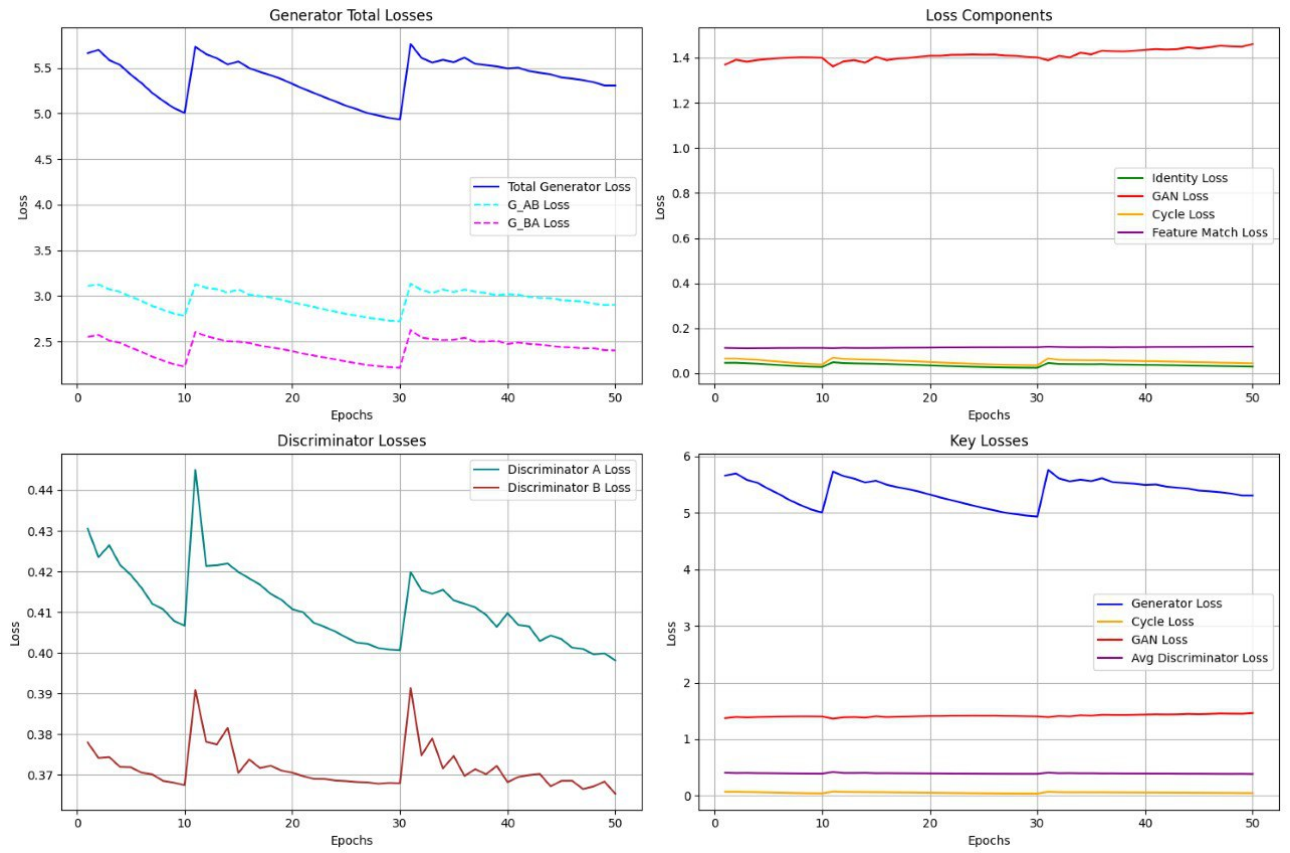


Figure 1.8: Loss curve for Stage 2 training

2. Task 2

2.1 Introduction

In this project, we explore the task of translating the appearance of real-world animals into the stylistic domain of Pokémon sprites. Our approach is based on a similar architecture previously used in Task 1, adapting it to handle the challenges of domain transfer between realistic and heavily stylized imagery. We will also dive into other experiments where we have encountered failures.

2.1.1 Configuration

All optimisers used Adam with a Cosine Annealing Warm Restarts scheduler. The following hyperparameters and loss weights were used for our best results:

- Random seed: 42
- Torch seed: 7193048279908058341
- Epochs: 501
- Batch size: 32
- Learning rate: $3e-4$
- Betas: (0.5, 0.999)
- $Scheduler_{t=0}$: 10
- $Scheduler_{t=mult}$: 2
- $Scheduler_{min-lr}$: $1e-5$
- λ_{cyc} : 6.0
- λ_{id} : 3.0
- λ_{gan} : 1.0
- λ_{gp} : 10.0
- λ_{fm} : 8.0
- $\lambda_{perceptual}$: 2.0

2.2 Architecture

2.2.1 Generator

We used a similar architecture as the one used in Task 1 (1.2.1) for our best results. The following parameters were also used for our generator:

- Residual Blocks: 6
- Use Attention: True
- Features Dimension: 64

Failure From VQGAN Implementation: **Vector Quantization Block**2.6 is used to convert the input image into a discrete representation, which is then passed through a series of residual blocks. The output of the residual blocks is then passed through a decoder to reconstruct the image. The vector quantization block is used to learn a set of discrete codes that can be used to represent the input image, allowing for more efficient encoding and decoding.

2.2.2 Discriminator

We used the same general architecture as the one used in Task 1 (1.2.2), with a few modifications specific to catering for the smaller images. Since our Pokémon sprite sizes are only 96x96, we removed the last down-sampling layer of the discriminator so bottleneck/features will be capped at $[B, 256, 12, 12]$. The reasons are two-fold:

- The structure of an animal and pokemon being very different, so we will not want to make the discriminator too powerful.
- A $[B, 256, 6, 6]$ feature map will not provide enough useful information to the generator either, thus a $[B, 256, 12, 12]$ feature map can better help the generator to identify the Pokémon sprite style.

2.3 Data Preparation

2.3.1 Collection

Due to our task of transferring the styles from a real-life animal to that of a Pokémon sprite, we were unable to find a paired dataset between these 2 domains. Thus, we resorted to utilising 2 individual datasets: one for the animal images, another for the Pokémon sprites, and then combining them together.

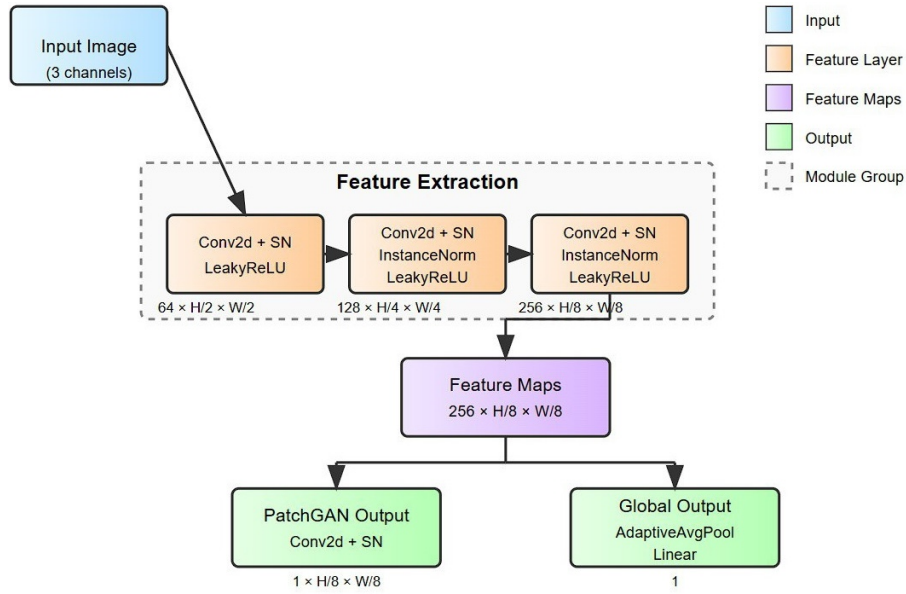


Figure 2.1: Task 2: Discriminator Architecture

Animal Dataset

The animal dataset ¹ we used consists of 5400 images of animals from 90 different classes / types of animals. We decided to use all images from the animal dataset. (Usability: 10.0 Total Image: 5.4k)

Pokémon Sprite Dataset

We decided to use Pokémon sprites as compared to the original Pokémon art as there is a distinct style to how Pokémon sprites are. Furthermore, using the original Pokémon art would not be suitable for a "style transfer" as we felt that the original art is at the prerogative of the individual artists when they created them. Thus, the styles across Pokémon will not be consistent, and our model will find it very difficult to learn the mapping between the two domains. Our failed attempts can be seen in Figure 2.2 using this dataset²



Figure 2.2: Unacceptable style transfer with Pokémon Official Art

Notice that the generated image does not resemble Pokémon art at all. We suspect that it is trying to colour the animal images in but it does not know exactly in what colour to use, and it does not make any informed or uninformed strong decision on how to colour. Thus, it ended up looking like a "mash-up" of the two domains, using the exact same colouring pattern.

¹Animal Image Dataset (90 Different Animals) – Zooming in on Wildlife: 5400 Animal Images Across 90 Diverse Classes <https://www.kaggle.com/datasets/iamsouravbanerjee/animal-image-dataset-90-different-animals>

²All Pokémon Official Artworks – A Comprehensive Dataset of All 1010 Official Pokémon Artworks <https://www.kaggle.com/datasets/hesselaar/all-pokemon-official-images>

Thus, we finalised on using the Pokémon sprite dataset ³ that consists of many different variations of the same Pokémon, namely those front-facing, back facing, shiny and normal. We have decided to use the normal and front-facing versions as those are the most suitable given that our animal datasets also exhibit the same characteristics and artstyle. (Usability 5.29, Total Image: 10.2k)

2.3.2 Pre-Processing

Due to the domains being vastly different, and the images are of different formats, we had to pre-process the images before feeding them into the model so that it will learn what it needs to learn.

Background Removal

The Pokémon sprites from the dataset are placed on a white background, whereas the animal images are not. Thus, we need to remove the background of the animal images so that the model can focus on the animals themselves rather than including how to convert the background to white.

To do so, we used a pretrained model (VGG16) to remove the background of the animal and convert them into white background found in `clean_background.py`.

No Tightly Paired Images

We **randomly** paired each animal to a **random** Pokémon and we ensure that each Pokémon is not used too many times, as the different Pokémons are not equally represented in quantity. Since we have 1025 unique Pokémon, with different quantity of front-facing Pokémons, and 5400 different animal images, we try to ensure that each type of Pokémon sprite is at most represented 6 times. Finally, we have a total of 5400 animal-Pokémon paired images.

Failure From Tight Pairs: Before choosing to adopt a more "unsupervised learning" approach, we tried to mimic a "supervised learning" by mapping the animals to a specific Pokémon based on how similar their general designs are. However, we found that the results were **not satisfactory**. Although it does look like a Pokémon sprite from afar, it has a much worse GMS score as compared to the non-tightly paired dataset. Furthermore, we also conducted an anonymised field test through telegram polling with students across various schools-NUS, SMU, NTU-where some specifically are involved in Pokémon clubs, and across various undergraduate years. As of writing, we have 131 responses, with 1.5x more people prefer the results from the randomly-paired dataset compared to the tightly paired dataset. The methodology of the poll can be found in Appendix 2.5. Thus, it was sufficiently convincing for us to use the randomly paired dataset method.

2.3.3 Data Split

We split the pre-processed dataset into a simple 90/10 train-test split. So, we will have 4860 train images and 540 test images. The training set (4860 images) is split further to a 99/1 split, where 4811 images will be used for training, while 49 will be used for validation (visualisation of how the pokemon will look like).

2.4 Loss Functions

Our loss functions for task 2 is identical to that of task 1 **with the exception of an additional loss, Perceptual Loss 2.4.2, for our Generator**. In summary, we used a combination of loss functions beyond the standard set of losses in a CycleGAN.

2.4.1 Discriminator Losses

We used a combination of loss functions beyond the standard adversarial loss to improve the performance of our discriminator. Our overall loss function for the discriminator is as follows:

$$L_D = L_{GAN} + \lambda_{gp} L_{gp} \quad (2.1)$$

³Pokemon sprite images – Total 10,437 Pokemon sprite images in 96x96 resolution.
<https://www.kaggle.com/datasets/yehongjiang/pokemon-sprites-images>

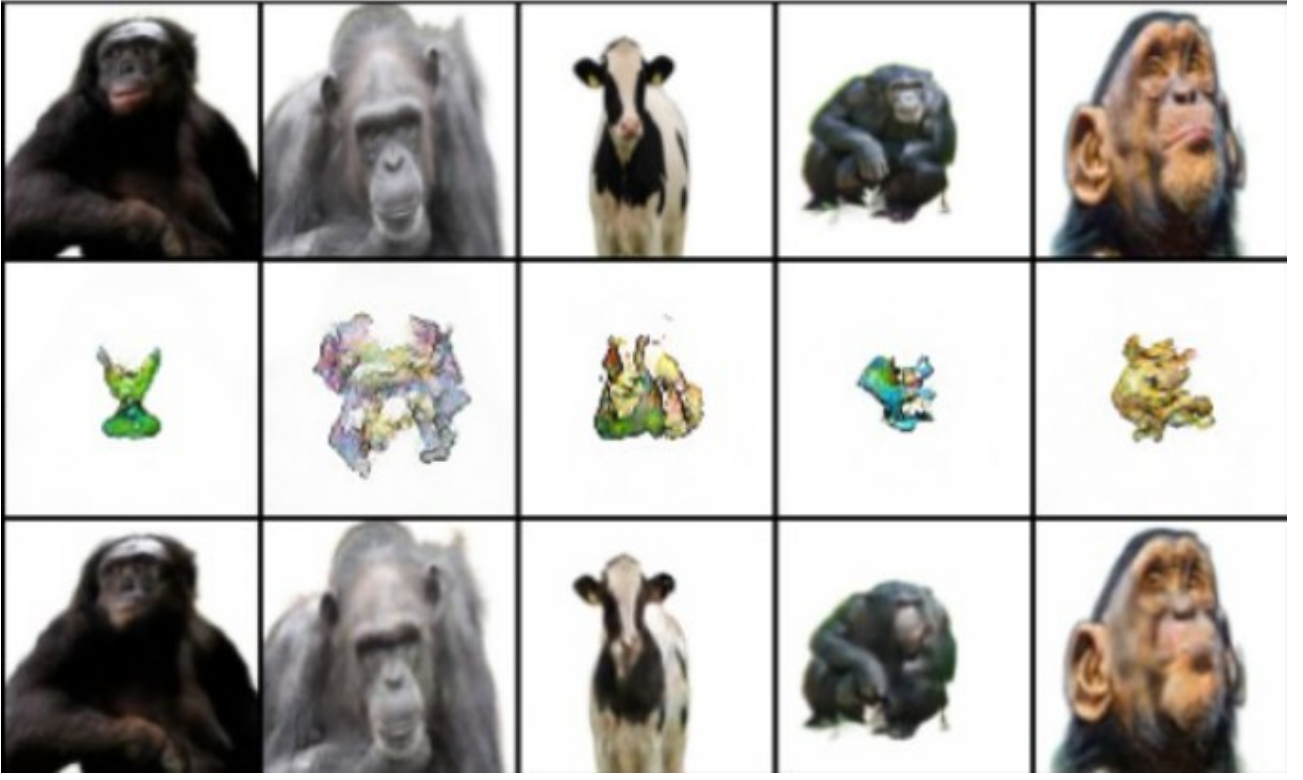


Figure 2.3: Unsatisfactory results from tightly paired dataset

Gradient Penalty – L_{gp}

We also use a gradient penalty, inspired from Wasserstein GAN (WGAN)⁴ in order to prevent the discriminator from learning too fast. Doing so avoids problems like mode collapse with vanishing gradients, and ensures that the generator is always learning, even if it is not performing well enough to "beat" the discriminator yet, thus improving stability of training.

$$\text{Gradient Penalty Loss} = 0.7 \cdot \mathbb{E} \left[(\|\nabla_{\hat{x}} D_{\text{patch}}(\hat{x})\|_2 - 1)^2 \right] + 0.3 \cdot \mathbb{E} \left[(\|\nabla_{\hat{x}} D_{\text{global}}(\hat{x})\|_2 - 1)^2 \right]$$

Relativistic – L_{GAN}

We extend the typical adversarial loss by using a relativistic loss, which compares the logits of discriminating the real and fake images in a relative manner on a global scope with Binary Cross Entropy Loss. Furthermore, we took inspiration from PatchGAN⁵ which whether the images is real or fake based on smaller patches, a local scope. Thus, we have another source of information for determining real or fake, which is computed with MSE loss. The equation for our relativistic loss for our "real" side is as follows, and the opposite will be applied for the "fake" side i.e. $D_{\text{real}} \leq D_{\text{fake}}, \text{True} \leq \text{False}$:

$$L_{GAN} = 0.7 \cdot \text{MSE}(D_{\text{real}} - \bar{D}_{\text{fake}}, \text{True})_{\text{patch}} + 0.3 \cdot \text{BCE}(D_{\text{real}} - \bar{D}_{\text{fake}}, \text{True})_{\text{global}} \quad (2.2)$$

2.4.2 Generator Losses

The generator also makes use of the **Relativistic** loss, which is an extension of the standard adversarial GAN loss. It also takes advantage of a pre-trained VGG16 model to compare the similarities of extracted, as referenced earlier as **Perceptual Loss**. Our overall loss function for the generator is as follows:

$$L_G = L_{GAN} + \lambda_{\text{perceptual}} L_{\text{perceptual}} + \lambda_{\text{cyc}} L_{\text{cyc}} + \lambda_{\text{id}} L_{\text{id}} + \lambda_{\text{feat}} L_{\text{feat}} \quad (2.3)$$

⁴Arjovsky, M., Chintala, S., & Bottou, L. (2017, July 17). Wasserstein Generative Adversarial networks. PMLR. <https://proceedings.mlr.press/v70/arjovsky17a.html>

⁵Isola, P., Zhu, J., Zhou, T., & Efros, A. A. (2016). Image-to-Image Translation with Conditional Adversarial Networks. arXiv (Cornell University). <https://doi.org/10.48550/arxiv.1611.07004>

Perceptual Loss – $L_{perceptual}$

We used a perceptual loss, which is a combination of the content and style loss. The content loss is computed using the features extracted from the VGG16 model, while the style loss is computed using the Gram matrix of the features. The perceptual loss helps to ensure that the generated images are perceptually similar to the target images, while also maintaining the overall structure and content of the images.

Cycle Consistency – L_{cyc}

As we would like the fake images to be able revertible to its original image i.e. the pixels are identical, we maintain using L1 Loss for our Cycle Consistency.

Identity – L_{id}

Similarly, we also use L1 Loss for our Identity loss since images in the same domain should not be altered.

Feature – L_{feat}

To encourage perceptual similarity between generated Pokémon and their corresponding animal inputs, we incorporate a Feature Loss (L_{feat}) based on intermediate activations from our discriminator. While perceptual loss is often computed using pre-trained networks like VGG, these models are trained on datasets such as ImageNet, which do not share the stylistic or semantic characteristics of our target domain. Instead, we extract features from our task-specific discriminator, which learns domain-relevant representations during adversarial training. We use the L1 norm to compute the distance between features of the generated and target images to focus on accurate reconstruction, as it is more robust to outliers and tends to preserve texture and structure better than L2, which can overly penalize large deviations and lead to smoother, less detailed outputs.

2.5 Methodology

Similar to task 1, we also made use of **Image Pooling** (as seen in 1.4.2). However, we did not employ any additional special training techniques due to the lack of time to experiment. Instead, we opted for the standard training procedure with the hyperparameters as mentioned in 2.1.1.

2.6 Metrics

Table 2.1: Evaluation Metrics for Animal \rightarrow Pokémon Style Transfer

Model	FID ↓	IS ↑	GMS ↓	SSIM ↑	Perceptual ↓	Diversity ↑
Animal \rightarrow Pokémon	96.37582	2.50949 \pm 0.24199	6.19713	0.7262	3.0451	0.1451
Pokémon \rightarrow Animal	205.14336	3.26938 \pm 0.30883	7.92128	0.5977	2.7445	0.2408

2.7 Conclusion

In summary, our work inherits the learnings from Task 1 by introducing modifications tailored to the more challenging domain translation between animal images and Pokémon sprites. Key enhancements include the incorporation of a Vector Quantization Block for efficient encoding (see 2.2.1) and adjustments to the discriminator architecture to accommodate smaller image resolutions (see 2.2.2). Our data preparation strategy—featuring background removal and randomized pairing (see 2.3.2)—ensured the model could efficiently learn the differences between the two domains. Finally, by integrating an additional Perceptual Loss into our loss functions (see 2.4.2), we further improved the perceptual fidelity of the generated images while maintaining training stability as seen from our metrics, utilising both subjective and objective analysis (see 2.6). These innovations together present a robust solution for unsupervised style transfer between highly disparate visual domains.

2.8 Appendix

2.8.1 Randomly Paired vs Supervised Paired Dataset Preference

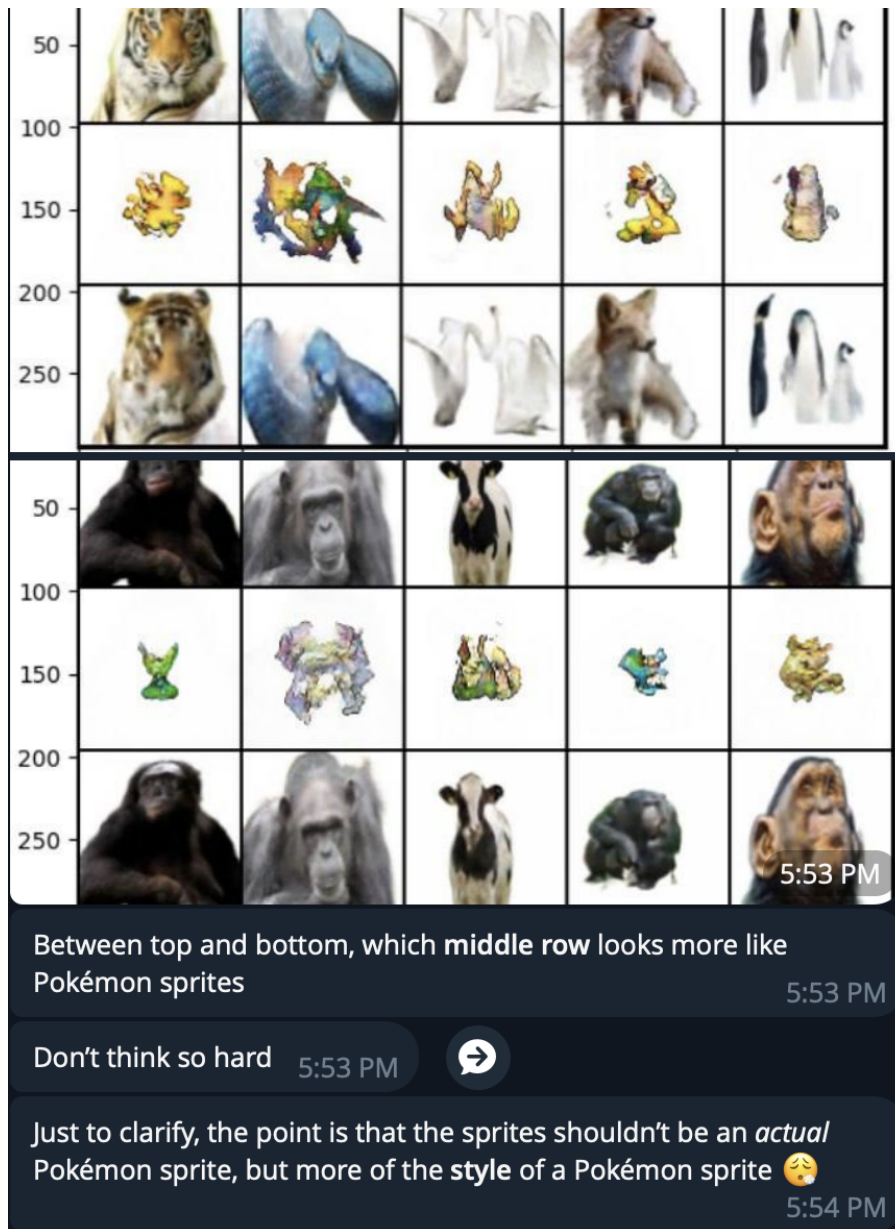


Figure 2.4: Context given to polling participants



Figure 2.5: Poll Results

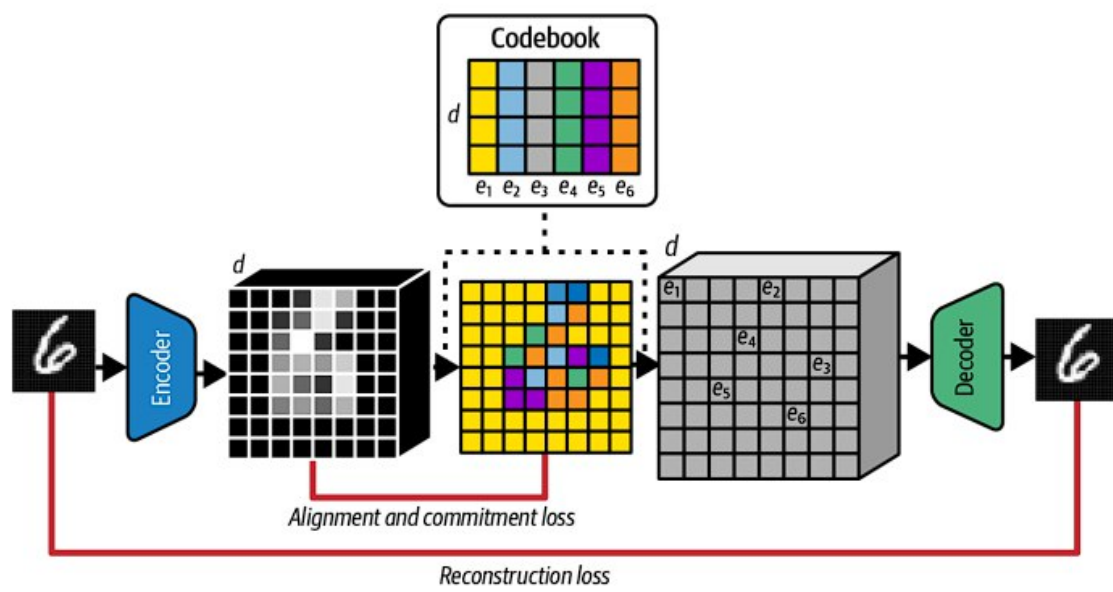


Figure 2.6: Vector Quantization Block

Cycle A: Real A \rightarrow Fake B \rightarrow Reconstructed A

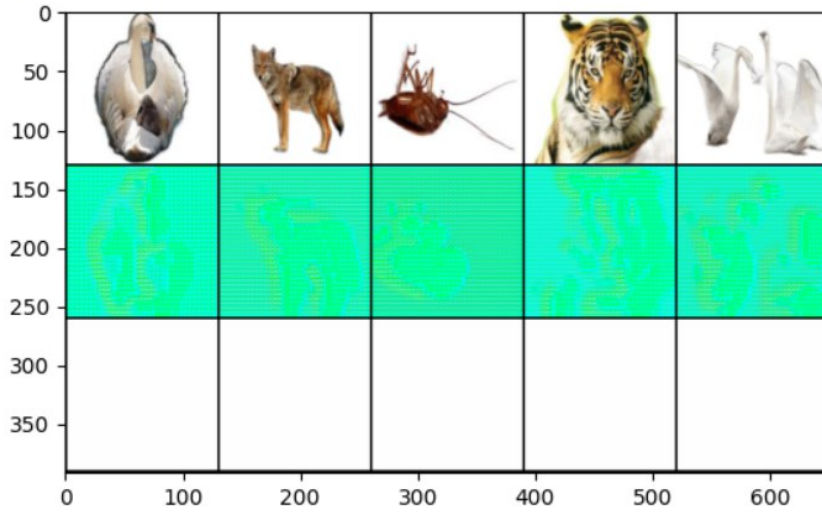


Figure 2.7: Complete failure of VQGAN for Animals \rightarrow Pokémon \rightarrow Reconstructed Animals

Cycle B: Real B \rightarrow Fake A \rightarrow Reconstructed B

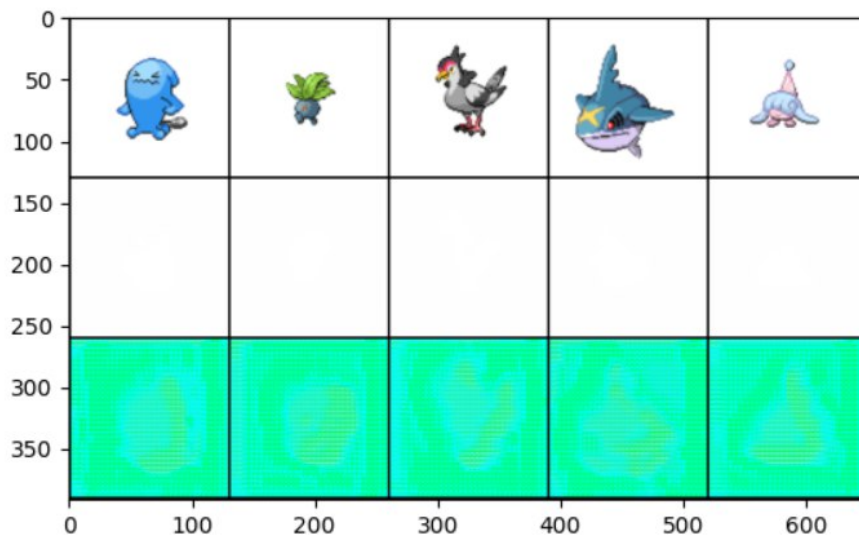


Figure 2.8: Complete failure of VQGAN for Pokémon \rightarrow Animals \rightarrow Reconstructed Pokémon