

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

Exploring CycleGAN and potential improvements for Unpaired Image-to-Image Translation

Karan Pardasani

kp955@scarletmail.rutgers.edu

Rutgers University

Department of Computer Science

Jaini Patel

jp1891@scarletmail.rutgers.edu

Rutgers University

Department of Computer Science

Animesh Sharma

animesh.sharma@rutgers.edu

Rutgers University

Department of Computer Science

Harish Udhay

harish.udhayakumar@rutgers.edu

Rutgers University

Department of Computer Science

Abstract

The class of algorithms in computer vision where the aim is to learn the mapping between an input image and output image is called image to image translation. There are two class of algorithms - paired and unpaired image to image translation. The CycleGAN architecture is used to train the model to learn the mapping using unpaired image to image translation. The goal for this training is to learn the mapping $G : X \rightarrow Y$ such that the transformation of the input image, $G(X)$ are very close to the images present in the target domain. Quantitative and Qualitative results are presented to evaluate the performance of the CycleGAN to learn the transformation from real to prerecorded synthetic images and vice versa, and also from real to live synthetic images and vice versa. After completing the training on standard CycleGAN model, we introduce some enhancements in the CycleGAN model. We train the enhanced model on same pairs of datasets that we train in case of standard CycleGAN model.

1. Introduction

The goal of Image-to-image translation (I2I) is to transfer images from a one domain to another while preserving the image representations. I2I has grabbed increasing attention and made massive progress in recent years due to its wide range of applications in many problems, such as image synthesis, segmentation, style transfer, restoration, and pose estimation. In this project, we experiment and present an overview of the I2I works developed in recent years. We will analyze the existing key machine learning frame works

(GANs). Additionally, we will experiment novel approach to enhance existing techniques.

2. Related Work

2.1. Recollection: Task 1

In Task 1, we implemented the basic GAN model with Generator architecture as described in [3] and implemented 70x70 PatchGAN for the Discriminator architecture as described in [2]. The Discriminator identifies whether the generated image belongs to the target distribution. The Generator uses Adversarial Loss to incorporate the effect of Discriminator on Generator weights. For this task, we performed the Qualitative Evaluation and Quantitative Evaluation by calculating Frechet Inception Distance [1, 5] and Inception Score [4]. Fig. 1 displays the output images generated from the model using real pizza dataset. From the images, we can conclude that in the generation of real images, the GAN is able to generate the correct shape and place the topping in the good locations. Also, the size of the toppings are good in relative to the size of the pizza. The model is not able to perform well in synthetic pizza generation. The shape is not proper and also the number of toppings are not consistent with the input image. Fig. 2 displays the output images generated from the model using pre-recorded synthetic pizza datasets.



Figure 1. Sample generated real pizza images by GAN



Figure 2. Sample generated synthetic pizza images by GAN

2.2. Recollection: Task 2

In Task 2, we focused on paired image-to-image translation and trained Pix2PixGAN model as described in [2] on Dayton dataset. The Generator uses Adversarial Loss to incorporate the effect of PatchGAN Discriminator on Generator weights. For this task, we performed the Qualitative Evaluation and Quantitative Evaluation by calculating IS and FID scores. Fig. 3 and Fig. 4 display the generated output images for overhead to street view translation and street view to overhead translation, respectively. We can conclude that in the generation of street view images, the GAN is able to generate different shapes of house, trees, roads and pavements. The visual acuity can distinguish between different objects but the distinction is not very clear. The model also has decent performance in generating overhead images. The model is able to correctly generate trees, roads and buildings. But, the distinction of different objects is not as apparent as the distinction in generated street view images.



Figure 3. Sample generated images of Street view



Figure 4. Sample generated images of Overhead view

3. Dataset

For the purpose of training CycleGAN and enhanced CycleGAN models, we have used datasets from three domains; pre-recorded synthetic pizza dataset, real pizza dataset and live synthetic pizza dataset. Our goal with unsupervised image-to-image translation is to translate the images from one domain to another, and so, we perform two types of translations; real to synthetic and synthetic to real, using two different synthetic datasets and one real pizza dataset. The total images in real pizza dataset train-test-val sets are 6688-836-824. The total images in pre-recorded synthetic dataset train-test-val sets are 13321-1538-1481.

To generate the live synthetic dataset, we have used 13 different types of toppings with each having approximately 5 to 6 images, 5 different types of pizza base images, and 6 different background images. We have applied alpha channel to the toppings and base images. Total number of images in the live synthetic dataset train-test-val sets are 13321-1538-1481.

We have pre-processed the images by up-scaling the image, then using random crop and finally normalizing it, for the train dataset. The final dimensions of all input images are 256×256 pixels.



Figure 5. Example of components of Live Synthetic images



Figure 6. Sample images from live synthetic pizza dataset



Figure 7. Sample images from real pizza dataset



Figure 8. Sample images from prerecorded Synthetic pizza dataset

4. Cycle GAN

We have implemented unpaired image-to-image translation using CycleGAN model as in the paper [6], on pizza datasets.

4.1. Methodology

The aim of CycleGAN is to learn a mapping $G : X \rightarrow Y$ as such the distribution of images from $G(X)$ is indistin-

guishable from the distribution Y using an adversarial loss. As this mapping is highly under-constrained, we couple it with an inverse mapping $F : Y \rightarrow X$ and establish a cycle consistency loss to push $F(G(X)) \approx X$ (and vice versa).

4.2. Formulation

Adversarial Loss:

This loss is used to learn the distribution of the target domain so that the Generator can generate images from source domain to the target domain.

$$\text{Loss_advers}(G, D_y, X, Y) = \frac{1}{m} \sum (1 - D_y(G(x)))^2$$

$$\text{Loss_advers}(F, D_x, X, Y) = \frac{1}{m} \sum (1 - D_x(F(y)))^2$$

Cycle Consistency loss:

This loss is used to make consistent information when we translate image from the source domain to the target domain.

$$\text{Loss_cyc}(G, F, X, Y) = \frac{1}{m} [(F(G(x_i)) - x_i) + (G(F(y_i)) - y_i)]$$

Identity loss:

This loss function is used to preserve information for the images that are already present in the target domain.

$$\begin{aligned} \text{Loss_identity}(G, F) &= E_{y \sim p_{data}}(y)[\|F(y) - y\|_1] \\ &+ E_{x \sim p_{data}}(x)[\|F(x) - x\|_1] \end{aligned}$$

Overall Loss:

$$\begin{aligned} \text{Overall Loss} &= \text{Loss_advers}(G, D_y, X, Y) \\ &+ \text{Loss_advers}(F, D_x, X, Y) \\ &+ \text{Loss_cyc}(G, F, X, Y) \\ &+ \text{Loss_identity}(G, F) \end{aligned}$$

4.3. Architecture

The architecture of CycleGAN has 2 Generator networks (G and F) and two discriminator networks, (D_X and D_Y). (fig. 9)

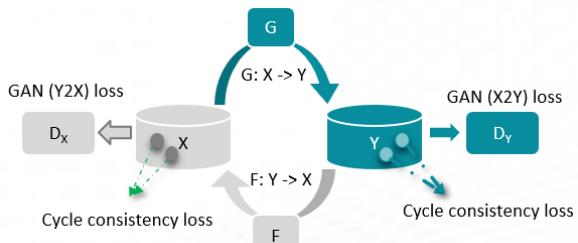


Figure 9. Cycle GAN architecture

Generator: c7s1-64, d128, d256, R256, R256, R256,

324 R256, R256, R256, R256, R256, u128, u64, c7s1-3.
 325 Also, for the discriminator, we use 70×70 patchGAN dis-
 326 criminator and the architecture is as follows:
 327 C64-C128-C256-C512

329 4.4. Training

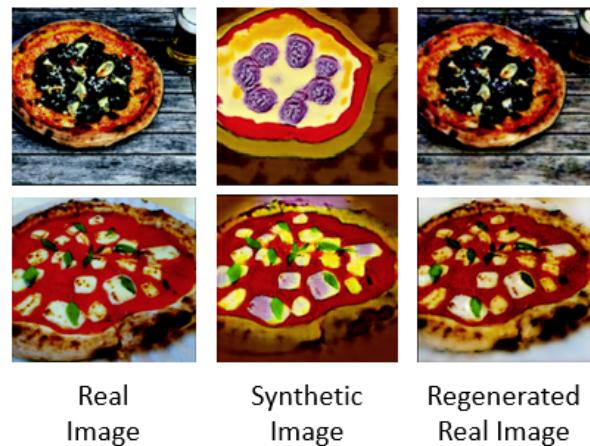
330 The CycleGAN was trained for two combinations of
 331 datasets as follows:

- 333 1. Live synthetic pizza dataset as one Domain and real
 334 pizza as the another. The live synthetic pizza dataset
 335 corresponds to images that were generated with top-
 336 pings, pizza bases and background that were manually
 337 collected.
- 338 2. Pre-recorded Synthetic pizza dataset as first Domain
 339 and real pizza dataset as another.

340 We have trained the model for 50 epochs and the FID scores
 341 are shown in the Fig 14, 15, 16 and 17.

342 4.5. Qualitative Evaluation

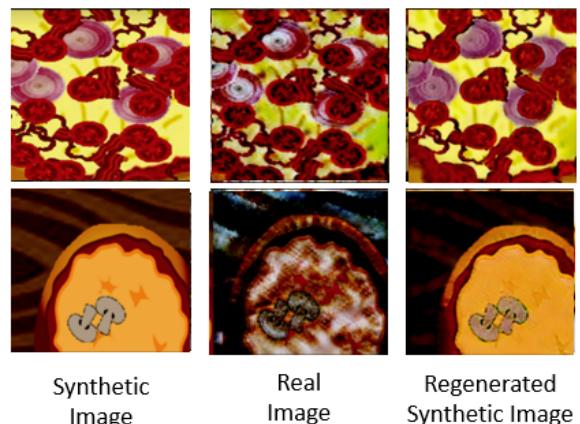
343 We perform Qualitative Evaluation of the models by
 344 manually checking the quality of the generated images.
 345 Here are few of the generated images for both the train-
 346 ing combinations as mentioned above. As we can see from
 347 Fig. 10, 11, 12 and 13, the quality of translated images
 348 from one domain to another is decent for all the combi-
 349 nations. There are some visible differences in the background
 350 texture, toppings texture, size and shape of the image ob-
 351 jects in the translated images from the original images. The
 352 reconstructed images are also similar to the original input
 353 images, but with some difference in proportion of the top-
 354 pings and pizza sizes.



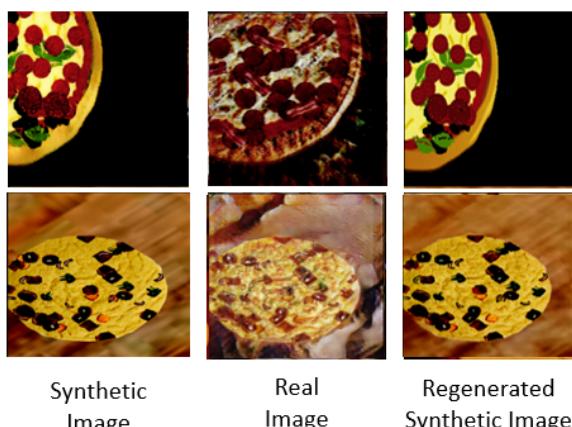
371 Figure 10. Translated Real images to synthetic pizza images and
 372 reconstructed it back to real Pizza from pre-recorded dataset

374 4.6. Quantitative Evaluation

376 For quantitative evaluation we use the following evalua-
 377 tion metric:



392 Figure 11. Translated Synthetic images to Real pizza images and
 393 reconstructed it back to Synthetic Pizza from pre-recorded dataset



409 Figure 12. Translated Syhthetic images to real pizza images and
 410 reconstructed it back to Syhthetic Pizza from live dataset



422 Figure 13. Translated real images to synthetic pizza images and
 423 reconstructed it back to real Pizza from Live dataset

425 FID - The Frechet Inception Distance (FID) , is an eval-
 426 uation metric that is used for evaluating the quality of gener-
 427 ated images and specifically developed to evaluate the per-
 428 formance of generative adversarial networks. For the Base
 429 CycleGAN, we observed that there was decrease in FID
 430 scores with increase in number of epochs, but after certain

number of epochs, the FID levels increased a bit, showing the possible effects of mode collapse. Fig. 14, 15, 16 and 17 show the FID scores for the base CycleGAN model.

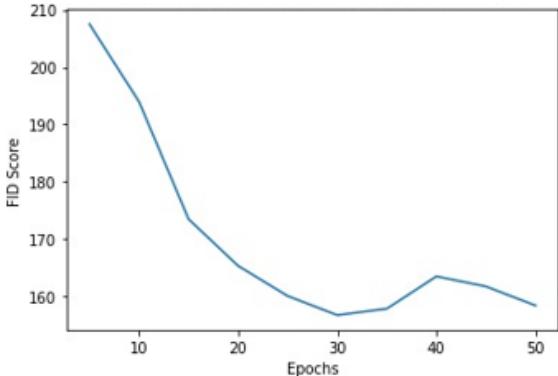


Figure 14. FID score for real images versus reconstructed real images from live pizza dataset

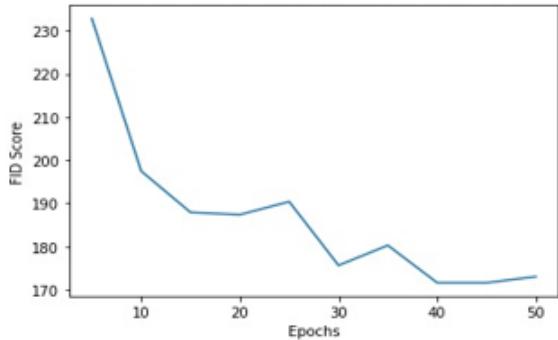


Figure 15. FID score for real images versus reconstructed real images from recorded pizza dataset

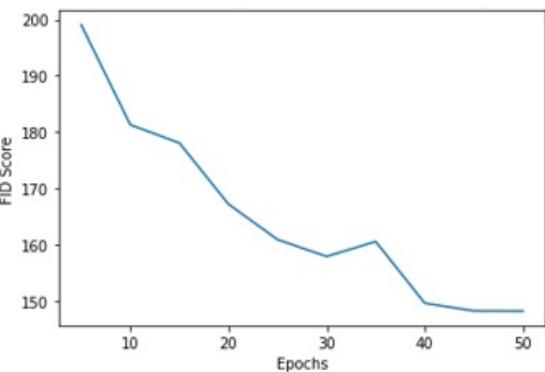


Figure 16. FID score for synthetic images versus reconstructed synthetic images from live pizza dataset

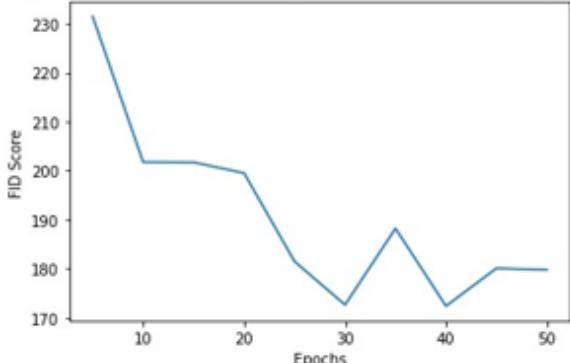


Figure 17. FID score for synthetic images versus reconstructed synthetic images from recorded pizza dataset

5. Improvements to Cycle GAN

5.1. Need for Improvements

While training the base Cycle GAN we noticed the following drawbacks:

- Cycle GAN does not produce good results if the difference between the two domains is too large. We believe this is partly due to the fact that the cycle consistency loss is very rigid and does a pixel-by-pixel comparison, which limits it's flexibility.
- The model suffered from mode collapse after some epochs. This happens when there is vanishing gradient in the discriminator. So, in this case, the weights of the discriminator stop updating and the generator keeps on getting better. So, generator will continue generating a specific type of images that will always fool the discriminator.

5.2. Introducing Global Discriminator

The standard architecture of Cycle GAN uses two PatchGAN discriminators. In addition to the two PatchGAN discriminators, we are introducing two Global Discriminators to the model. The PatchGAN Discriminator classifies whether each patch of the image belongs to the target domain distribution or not. The Global Discriminator classifies by looking at the whole image as a whole and outputs the probability of the image being in the distribution or not.

The image Fig. 18 shows the architecture of the enhanced CycleGAN model.

PatchGAN discriminator is used because it focuses on local features (finer details) and independently looks at local patches to decide whether the image is real or fake. Also, using only PatchGAN discriminator is insufficient since the patches of the images are dependant on each other. So, we are using the Global Discriminator to consider the relation

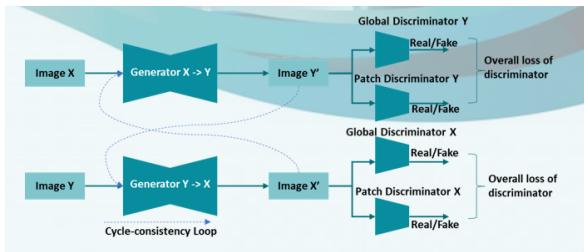


Figure 18. Cycle GAN architecture

between the patches of images. Also, PatchGAN is important because in considering the global features, the local information is sometimes overlooked. Hence, including both PatchGAN and Global Discriminator should allow the model to benefit from both types of Discriminators and make the model more robust against mode collapse.

5.3. Considering FID score in Generator Loss function

FID score is the similarity score used to measure the closeness of two datasets of images. Rather than directly comparing images pixel by pixel, the FID compares the mean and standard deviation of one of the deeper layers in the Inception v3 network. These layers are closer to output nodes that correspond to real-world objects and further from the shallow layers near the input image. As a result, they tend to mimic human perception of similarity in images. FID scores are inversely proportional to the similarity of the image sets. So, we consider FID score in loss function. Since the model tries to navigate towards minimising loss function, the model will have low values of FID. Now, there are three loss functions in the standard GAN model Adversarial Loss, Identity Loss and Cycle Consistency Loss and we add two more for the Enhanced GAN model:

- **Adversarial Loss:** This loss penalises if the Discriminator is able to differentiate between the generated images and the images from the target domain.
- **Identity Loss:** To calculate this loss, we take an image from the target domain and generate the image using the generator G . Since the input is in the target domain, the output image should ideally be the same as input image. So, this loss penalises if the generated image is very far from the input image, if input image is from the target domain. We consider L1 loss to penalise the divergence of input image from the target domain and the output image.
- **Cycle Consistency Loss:** To calculate this loss, we take the image from the source domain and generate the output image that should be in target domain. Also, we use this output image to reconstruct the input image using another Generator that maps the output domain

to input domain. The reconstructed image should be very close to output image. Hence, this loss function penalises if the input image is very far from the reconstructed image. We consider L1 loss to penalise the divergence of input image and the reconstructed image.

- **FID for Identity Loss:** The idea is that in addition to the L1 Loss from the Base CycleGAN model, using the FID between the input and output images will produce a more robust identity loss to penalise the divergence of the input image from the target domain and the output image. We scale the FID score down by a user-tunable parameter and then take the weighted mean of the L1 Loss and the FID score as the net identity loss.
- **FID for Cycle Consistency Loss:** Now, in the Cycle Consistency Loss, it was discussed that we use L1 loss to penalise the divergence of the input image from the source domain and the reconstructed image. In this case, we are considering the FID between the two images. Penalising the Generator on large FID scores between the input from source domain and the corresponding reconstructed image will make the images similar to each other.

The Final Loss function of the Enhanced Generator Model is:

$$\begin{aligned}
 \text{Updated Generator Loss} = & \text{L1 Identity Loss} \\
 & + \text{FID Identity Loss} \\
 & + \text{Adversarial Loss} \\
 & + \text{L1 Cycle Consistency Loss} \\
 & + \text{FID Cycle Consistency Loss} \quad (1)
 \end{aligned}$$

5.4. Generator Architecture

We adopted the Generator architecture from the Baseline CycleGAN model. We have two generators $G : X \rightarrow Y$ and $F : Y \rightarrow X$. The loss function of the generator is described in equation (1). The weights of the two generators are updated according to the equation (1).

5.5. Discriminator Architecture

For the discriminator network two 70×70 PatchGANs were used, which aim to classify whether 70×70 overlapping image patches are real or fake. Such a patch-level discriminator architecture has fewer parameters than a full-image discriminator and can work on arbitrarily sized images in a fully convolutional fashion. We have adopted the PatchGAN Discriminator Architecture as described in [2].

Likewise, we also added two global discriminators - classifying if entire image is fake or real. The global and patch

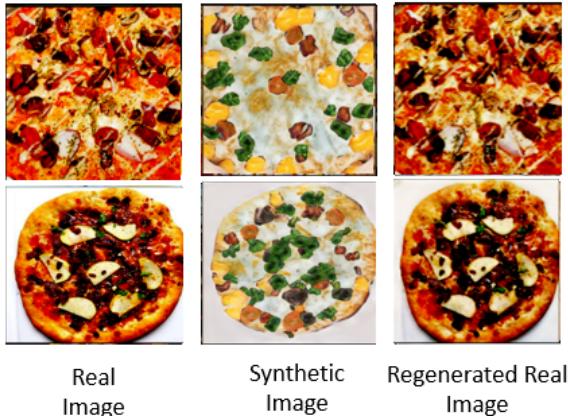
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

648 discriminator work independently and the net Generator
 649 Adversarial Loss includes components from all the discrimi-
 650 nators to make the model robust (or immune) to mode col-
 651 lapsed.
 652

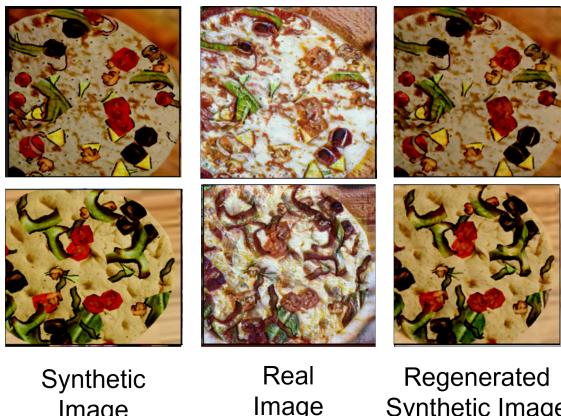
653 5.6. Qualitative Evaluation

654 For the training of Enhanced model, we use the batch
 655 size of 8.

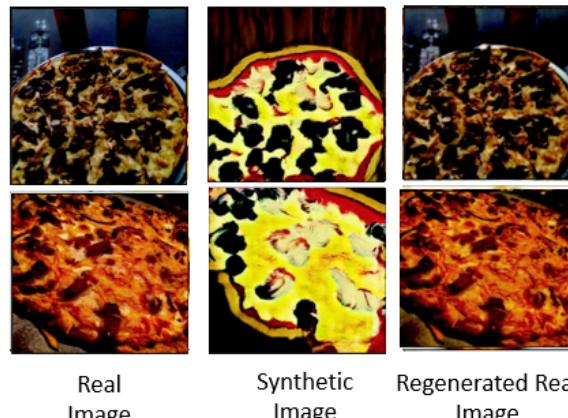
656 We perform Qualitative Evaluation of the enhanced
 657 model by manually checking the quality of the generated
 658 images. Here are few of the generated images for both the
 659 training combinations as mentioned above. As we can see
 660 from Fig. 19, 20, 21 and 22, the quality of translated im-
 661 ages from one domain to another is better than from the
 662 base CycleGAN model for all the combinations. The new
 663 model is able to capture the proportions, size, texture and
 664 shape of toppings of well compared to the Base model in
 665 the synthetic as well as reconstructed images.



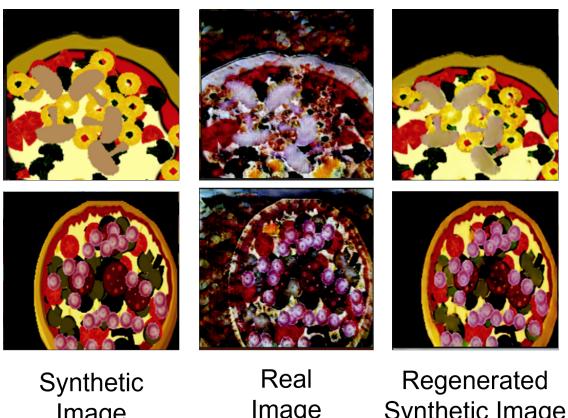
666 Figure 19. Sample output images of live dataset from Enhanced
 667 Cycle GAN. Real → generated Synthetic → reconstructed Real



697 Figure 20. Sample output images of live dataset from Enhanced
 698 Cycle. Synthetic → generated Real → reconstructed Synthetic



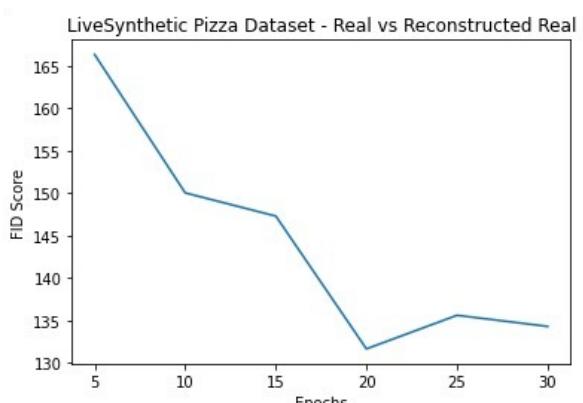
702 Figure 21. Sample output images of pre-recorded dataset from En-
 703 hanced Cycle. Real → generated Synthetic → reconstructed Real



704 Figure 22. Sample output images of pre-recorded dataset from En-
 705 hanced Cycle. Real → generated Real → reconstructed Synthetic

706 5.7. Quantitative Evaluation

707 Below graphs show FID scores from the enhanced model
 708 evaluation. (Fig. 23, 24, 25 and 26).



709 Figure 23. Enhanced CycleGAN - FID Score for Real images ver-
 710 sus reconstructed real images for Live Synthetic pizza dataset

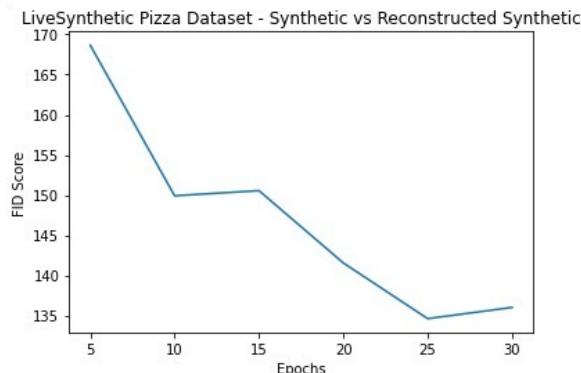


Figure 24. Enhanced CycleGAN - FID Score for synthetic images versus reconstructed synthetic images Live Synthetic pizza dataset

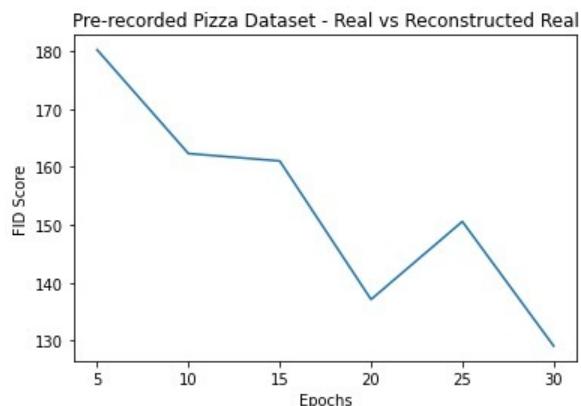


Figure 25. Enhanced CycleGAN - FID Score for real images versus reconstructed real images for pre-recorded pizza dataset

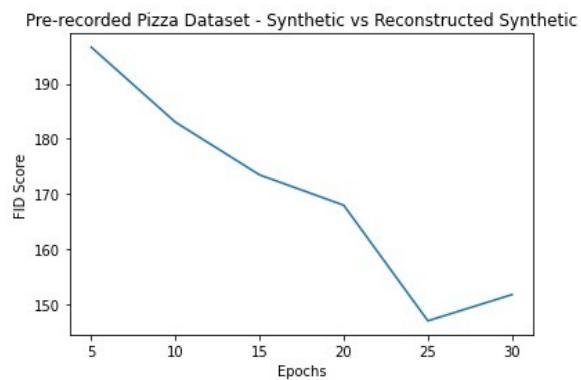


Figure 26. Enhanced CycleGAN - FID Score for synthetic images versus reconstructed synthetic images for pre-recorded pizza dataset

From the above visualizations of FID scores of two datasets, we infer that the range of FID scores for Live dataset are lower, compared to the pre-recorded dataset. Hence, the model is better with live dataset than synthetic dataset.

6. Discussion and Limitations

With our enhancements, we observed that the identity and Cyclic Consistency losses are more robust when FID is included, especially for larger batch sizes. Also, the addition of the Global Discriminator helps make the model more immune to mode collapse resulting from stagnation of learning in discriminator. So, overall the quality of the images have improved. Although the model performed better than the Base CycleGAN model, it has slower training time due to additional calculations that have to be performed for FID scores calculations and also for training the two global discriminators.

7. Conclusion

To conclude, the results of this paper suggest that, although CycleGAN is considered as one of the powerful unpaired image-to-image translation techniques, it fails to perform well when the domain gap between the source and target domains is large. The enhanced model proposed in this paper, tries to overcome some of the limitations of the CycleGAN and achieves better results, both qualitatively and quantitatively but it takes increases the training time and also requires more resources than the baseline model(CycleGAN).

8. Contribution

The rough outline of how we worked in Task 3 and 4 is as follows:

1. Task 3: Training and Evaluation of Cycle GAN
 - (a) Pre Recorded to Real Pizza and vice versa: Animesh and Harish
 - (b) Real Pizza to Pre Recorded Pizza and vice versa: Karan and Jaini
2. Task 4: Enhancement of Cycle GAN
 - (a) Generation of Live Dataset: Jaini, Animesh, Karan and Harish
 - (b) Enhancement of FID: Harish and Animesh
 - (c) Global and PatchGAN Discriminator: Karan and Jaini
3. Presentation : Animesh, Harish, Karan and Jaini
4. Report: Animesh, Harish, Karan and Jaini

Everyone gave equal contribution and everyone is satisfied with each other's participation and support.

9. Code Repository and Supplementary Results

The code for our Step 3 and Step 4 can be found at: [GitHub Repository](#). We have added another file containing Inception Scores which further supports the better performance of our improved model.

756	810
757	811
758	812
759	813
760	814
761	815
762	816
763	817
764	818
765	819
766	820
767	821
768	822
769	823
770	824
771	825
772	826
773	827
774	828
775	829
776	830
777	831
778	832
779	833
780	834
781	835
782	836
783	837
784	838
785	839
786	840
787	841
788	842
789	843
790	844
791	845
792	846
793	847
794	848
795	849
796	850
797	851
798	852
799	853
800	854
801	855
802	856
803	857
804	858
805	859
806	860
807	861
808	862
809	863

864	References	918
865		919
866	[1] A. Borji. Pros and cons of GAN evaluation measures: 867 New developments. <i>CoRR</i> , abs/2103.09396, 2021. 1	920
868	[2] P. Isola, J. Zhu, T. Zhou, and A. A. Efros. Image- 869 to-image translation with conditional adversarial net- 870 works. <i>CoRR</i> , abs/1611.07004, 2016. 1, 2, 6	921
871	[3] J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses 872 for real-time style transfer and super-resolution. volume 873 abs/1603.08155, 2016. 1	922
874	[4] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, 875 A. Radford, and X. Chen. Improved techniques for 876 training gans. June 2016. 1	923
877	[5] M. Seitzer. pytorch-fid: FID Score for Py- 878 Torch. https://github.com/mseitzer/ 879 pytorch-fid, August 2020. Version 0.2.1. 1	924
880	[6] Y. Zhao, R. Wu, and H. Dong. Unpaired image-to- 881 image translation using adversarial consistency loss. 882 <i>CoRR</i> , abs/2003.04858, 2020. 3	925
883		926
884		927
885		928
886		929
887		930
888		931
889		932
890		933
891		934
892		935
893		936
894		937
895		938
896		939
897		940
898		941
899		942
900		943
901		944
902		945
903		946
904		947
905		948
906		949
907		950
908		951
909		952
910		953
911		954
912		955
913		956
914		957
915		958
916		959
917		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971