

000
001
002054
055
056

PokeDog CycleGAN- Experiments in Encouraging High Level Feature Change in CycleGAN Architecture

Author: Avidesh Marajh

003
004
005
006
007
008
009
010
011
012
013
014057
058
059
060
061
062
063
064
065
066
067
068

Anonymous CVPR submission

015

Paper ID *****

016

017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034069
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

Abstract

Cycle-Consistent Adversarial Networks(CycleGAN) are a variation of GAN architecture which effectively performs unpaired image translation. This architecture has seen great results in translating between classes which have low level feature differences such as styles or patterns. It has, however, seen great difficulty in high level feature changes, specifically geometric changes. In this paper, I shall outline attempts at improving the CycleGAN architecture in order to better achieve geometric changes. Using the problem of unpaired translation between images of dogs and images of dog-like entities from the popular media franchise, Pokemon, I shall demonstrate several proposed methods for a novel CycleGAN architecture which can achieve geometric changes. These methods includes additional terms for the loss functions of both the generators and discriminators as well as random grayscaling of input images. These methods achieve unsatisfactory results in geometric changes to the source image but do increase the fidelity of outputs.

035
036
037

1. Introduction

038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

Pokemon is an iconic part of many young adults' childhoods. Known for its adorable character designs and distinct art-style, people have interacted with this franchise through video games, anime, trading card games and various other merchandise. It comes as no surprise that it is the highest grossing media franchise of all time, leading its runner-up by 20 billion dollars. Be it nostalgic fans or eager pet parents, many people would certainly find joy in seeing their beloved pets in the style of these lovable Pokemon. This project, inspired by my dogs Onix and Riolu, seeks to achieve just that.

In this paper, I present my findings in attempting to apply Cycle-Consistent Adversarial Networks(CycleGAN) to this problem of unpaired image translation between dogs and Pokemon. Generative Adversarial Networks are generative

models which consists of a generator and a discriminator which learn simultaneously. The generator learns to generate outputs of a specified class and the discriminator learns to distinguish between fake generator outputs and real datapoints from the dataset. They are driven by an adversarial loss, such that each model minimizes its loss by outperforming the other.

The CycleGAN architecture is a method for unpaired image translation which utilizes two discriminators and two generators. For two classes of images, X and Y, GeneratorXY takes an input of a real image from class X and produces a fake image of class Y. Conversely, GeneratorYX takes a real input from Y and produces a fake image of X. The discriminatorX aims to discern the outputs from GeneratorYX from real datapoints of X and discriminatorY aims to discern the outputs from GeneratorXY from real datapoints of Y. What truly defines the CycleGAN architecture is its multi-term loss calculation. The discriminators are still driven by adversarial loss from their corresponding Generator, but the Generators are driven by, in addition to adversarial loss, a cycle-consistency loss and an identity loss. The identity loss penalizes the generators for a metric of difference between an output y' for an input image y of class Y to the generatorXY (ideally it would perform no change for an input y). The cycle-consistency loss penalizes the generators for a metric of difference between an input x of class X and an output x' , where x' is the output of the generatorYX for an input y' and where y' is the output of the generatorXY for the input x .

These loss terms define the unique performance of CycleGAN's which allow them to perform image translation tasks with little information loss. That is, the model preserves high level features and performs minimal change. Correspondingly, the CycleGAN architecture excels in translating low level features(shapes and patterns) and is unsuitable for geometric changes and high level feature changes(recognizable artifacts or large components). This, works against our task of dog to Pokemon conversion since,

108 while Pokemon are mostly based on real animals, they usually have distinct shapes and features. My reasoning for
 109 selecting the CycleGAN architecture is its affinity for unpaired image translation.
 110 Had I used a paired image translation method, such as styleGAN, it would simply imitate
 111 the style of a provided Pokemon input onto the input image
 112 of a dog. The use of unpaired image translation allows for
 113 the creation of a Pokemon which has modified but recognizable
 114 features of the inputted dog. This will allow for the creation of a more original seeming Pokemon which closely
 115 resembles the inputted dog.

116 As such, this project served as an experiment in methods for
 117 encouraging high level feature changes in CycleGAN architecture
 118 using a variety of original methods. Some of these
 119 methods include, grayscale transformation of input images
 120 at a decaying random rate, a variety of new loss terms for
 121 both the discriminator and generator, and alteration to the
 122 architecture which involves replacing the two discriminators
 123 with a single classifier.

124 While none of these methods showed significant improvements
 125 in feature changes, they did show higher fidelity outputs than the base CycleGAN architecture.

131 2. Background

132 2.1. Unpaired Image-To-Image Translation Using 133 Cycle-Consistent Adversarial Networks (Zhu 134 et al. 2223-2232)

135 This paper introduces the CycleGAN architecture which
 136 is employed in our model. Naturally our model faces many
 137 of the limitations addressed in this paper. Most notably, this
 138 architecture is shown to be ineffective for problems which
 139 require geometric changes. This does limit the fidelity of
 140 dog/Pokemon translations as the Pokemon are not always
 141 quadrupedal and may have various additional features not
 142 common to dogs.

143 While the author makes no suggestions for addressing this
 144 issue, I attempt to encourage large feature changes through
 145 the use of additional loss terms.

146 2.2. PokeGAN: P2P (Pet to Pokemon) Stylizer 147 (Hedge et al.)

148 This project aimed toward a very similar objective as this
 149 one. That is, it aimed to transfer the styles of Pokémon onto
 150 various pictures of animals. Similar to my own, their model
 151 applied blue and green tints on output images in order to
 152 generate ‘Pokémon’. This is explained by the authors as
 153 being due to these being the most commonly occurring colors
 154 in the anime. This is very similar to the issues encountered
 155 by our baseline model in which a green tint was being
 156 applied to outputted dog images. This issue was left unad-
 157 dressed by the authors and was dismissed as an outlier only
 158 for inputs with little color variation.

159 This is addressed in my project through the use of a random
 160 grayscale transformation with a decaying probability.

161 2.3. Unsupervised Pixel-Level Domain Adaptation 162 with Generative Adversarial Networks (Bous- 163 malis et al. 3722-3731)

164 This paper also addressed the issue of excessive augmentation
 165 of irrelevant feature in Style Transfer problems. This
 166 issue is apparent in our baseline model by the application
 167 of filters over the entire image, not only augmenting the
 168 dog/Pokémon but also greatly altering the background. The
 169 authors addressed this problem using a content-similarity
 170 loss which uses an image mask which excludes features
 171 which should not be changed and adding the mean squared
 172 error between real and generated images to the total loss.
 173 This was extremely effective but required a pre-existing
 174 mask which can isolate the features which are required to
 175 change, which is not possible for our problem as our input
 176 images feature entities in different poses.

177 In my project, I implemented a diff loss function which is
 178 meant to create a minimax objective which minimizes uniform
 179 change and introduces large, isolated changes.

180 2.4. From GAN to WGAN (Weng)

181 The defining feature of GAN is the adversarial simultaneous
 182 training of two different models, a generator and a discriminator. This can lead to issues, however, when
 183 the discriminator overtakes the generator in performance re-
 184 stricting the generator from attempting to learn effectively.
 185 In this paper, they discuss multiple solutions to pro-
 186 moting stable learning. In addressing the problem of dis-
 187 criminator dominance, they suggested softer labels for the
 188 discriminator loss, eg. 0,1,0.9 instead of 0,1.

189 In my project, I attempt to address this issue using a class
 190 loss, which provides a secondary learning objective for the
 191 discriminator, making it more difficult for it to converge.

192 2.5. Augmented CycleGAN: Learning Many-to- 193 Many Mappings from Unpaired Data (Alma- 194 hairi et al.)

195 This paper suggests a method for diverse outputs in
 196 GANs. Their method of introducing a noise vector as addi-
 197 tional input to the generators is very promising but I believe
 198 it is unsuitable for our problem which has a very diverse
 199 dataset and requires a focus on higher fidelity than diverse
 200 outputs.

201 For my project, I utilized random grayscaling in order to
 202 prevent the model from overfitting to the color mappings of
 203 the images in the dataset.

204 3. Dataset

205 For this project, I utilized two sources of data:

216

3.1. Dog dataset

217

For my dog images, I employed the Stanford dogs dataset. This dataset is built using images and annotations from ImageNet and is made for large machine learning projects such as this. This dataset contains 20,580 images of dogs from 120 categories, with a variety of poses and backgrounds.

224

<http://vision.stanford.edu/aditya86/>
ImageNetDogs

226

3.2. Pokemon dataset

227

For my Pokemon images, I scalped the anime screenshot of gallery from a Pokemon wiki site, Bulbapedia. Using a web request python script to iterate through their unorganized, but well labelled, gallery of Pokemon images, I downloaded all images which mentioned the name of Pokemon from a personally curated list of 50 dog-like Pokemon. This resulted in a dataset of 940 images. For this dataset, I relied on images from the Pokemon anime, rather than the video games or trading card games, due to the variety of poses and backgrounds which accompany screenshots from the anime. The code used for this can be found in the link below.

240

[https://github.com/Deshy23/Dog2Pokedog/](https://github.com/Deshy23/Dog2Pokedog/blob/main/scrape.py)
blob/main/scrape.py

242

3.3. limitations

243

By machine learning standards, my Pokemon dataset is very small and is greatly dwarfed by my dog dataset. This scarcity in available data greatly increase the difficulty of the task.

248

4. Methods

249

My networks are constructed using the following blocks: Feature Map Block: A single 2D convolutional layer with a kernel size of 7 and padding size of 3.

252

Contracting Block: A single 2D convolutional layer with a kernel size of 3, padding size of 1 and a stride of 2. This is followed by an instance norm and a ReLU activation.

255

Expanding Block: A single 2D convolutional layer with a kernel size of 3, padding size of 1 and a stride of 2. This is followed by an instance norm and a ReLU activation..

258

Residual Block: Two 2D convolutional layer with a kernel size of 3, padding size of 1. This is followed by an instance norm and a ReLU activation.

262

4.1. Baseline

263

For my baseline model, I used a standard implementation of the CycleGAN architecture, as follows:

266

Generators: Each Generator consisted of a feature map block which receives 3 input channels and 64 output channels, followed by 2 contracting blocks each of which doubles the

number of channels. This feeds into a series of 8 residual blocks each of which have 256 channels, input and output. This is followed by 2 expanding blocks, each of which halve the number of channels followed by a feature map block which takes an input of 64 channels and an output of 3 channels. This is finally subject to a Tanh activation function.

Discriminators: Each Discriminator consisted of a Feature Map Block with 3 input channels and 64 output channels, 3 Contracting Blocks, each of which double the number of channels and finally, a 2D convolutional layer which outputs 1 channel.

Each discriminator was optimized separately as defined by an adversarial loss function.

The generators were optimized together by a hybrid loss function of the adversarial loss (with weight 1.0), the identity loss (with weight 0.1) and the cycle consistency loss (with weight 10)

Adversarial losses used the Binary Cross Entropy as its criterion and Identity and Cycle losses used L1 loss.

This model was used as a baseline due to its verifiable performance when applied to the horse to zebra problem.

4.2. Diff Loss

My first modification of the architecture is an additional loss term for the Generator. In creating this loss term, my intention was to create a term which encourages the model to perform notable geometric changes. To do this I sought to disincentivise small uniform changes and encourage large localized changes. This was formulated as a quotient of the L1 loss and the Mean Squared Error as follows: For an inputted image of class X, x , to a GeneratorXY, and the corresponding output y' and where L_1 is the L1 loss and MSE is the Mean Squared Error,

$$L_{diff} = \frac{L_1(x, y')^2}{MSE(x, y')} \quad (1)$$

The objective of this function is to minimize the L1 loss and maximize the MSE loss between real inputted images and fake outputted images. This is due to the nature of these loss functions, where the L1 loss penalize changes uniformly but the MSE loss harshly penalizes having few pixels which have large differences. Therefore, this diff loss will be minimized by having localized changes in the image rather than uniform changes such as a tint. That is, the image should have few pixels which have large changes rather than many pixels with small changes.

4.3. Class Loss

One of the main problems observed to be faced was the imbalance in discriminator and generator performance. This is a common failure mode for GANs and to prevent the discriminator from reaching convergence so rapidly, an

324 additional term was added to its loss. This class loss, for a
 325 DiscriminatorX, takes an input of a real image of class Y, y,
 326 and a generated image of class X, x', where x' is the output
 327 of the GeneratorYX for an input y. It then finds the discrimi-
 328 nators prediction for each term and computes the L1 loss
 329 between these terms. That is,
 330

$$L_{class} = L_1(disc_X(x'), disc_x(y)) \quad (2)$$

331 This loss function seeks to minimize the difference in pre-
 332 diction for an image of class Y and its corresponding fake
 333 image of class X. This is to increase the complexity of the
 334 discriminator's task to prevent it from converging prema-
 335 turately while also encouraging it to base its decision on the
 336 identification of remnant features of class Y rather than ex-
 337 ploiting other things such as the image quality of the gener-
 338 ator outputs.
 339

340 4.4. Retention Loss

341 As an alternative to the Diff Loss, I also devised the re-
 342 tention loss, an additional loss function for the Generator.
 343 This loss function seeks to similarly encourage changes to
 344 geometric features by punishing the Generator for maintain-
 345 ing features from the input images. For this it leverages the
 346 Discriminator as follows:
 347

348 For a GeneratorXY with some output y' and the Discrimi-
 349 natorX where BCE is the binary cross entropy,
 350

$$L_{ret} = BCE(disc_X(y'), 0) \quad (3)$$

351 This loss function, using similar logic as the above Class
 352 Loss, encourages the Generator to ensure its output y' is
 353 not similar to real images from class X.
 354

355 4.5. Classifier

356 Thus far, my contributions have involved additional loss
 357 terms to shift the optimization goals of the model. This
 358 change however, makes a large architectural change. In-
 359 stead of utilizing 2 discriminators for each class, I utilize a
 360 single Classifier which aims to output 0 for images of class
 361 X and 1 for images of class Y. The logic behind this change
 362 is to have a Classifier which learns some set of underly-
 363 ing common features between the two classes and the dif-
 364 ferences between them and encourages the Generators to
 365 utilize this underlying mapping. This also eliminates the
 366 failure mode of having the discriminators converge before
 367 the generators and thus prematurely end learning, as the the
 368 classifier faces a more complex task than the individual dis-
 369 criminatoras.
 370

371 The Classifier uses the same architecture as a single Dis-
 372 criminator as described in the baseline model(3.1). The loss
 373 function for the new model are as follows:
 374

375 Generator adversarial loss: Let $val_X = 0$ and $val_Y = 1$:
 376

377 For x',y' where x',y' are generated images from Generato-
 378 rYX and GeneratorXY, and BCE is the Binary Cross En-
 379 tropy,
 380

$$\begin{aligned} L_{gen,adv} = & BCE(Classifier(x'), 1 - val_Y) \\ & + BCE(Classifier(y'), 1 - val_X) \end{aligned} \quad (4)$$

381 Classifier adversarial loss For x,y,x',y' where x,y are real
 382 images from X,Y and x',y' are generated images from Gen-
 383 eratorYX and GeneratorXY, and BCE is the Binary Cross
 384 Entropy,
 385

$$\begin{aligned} L_{classifier,adv} = & BCE(Classifier(x), val_X) \\ & + BCE(Classifier(x'), val_Y) \\ & + BCE(Classifier(y), val_Y) \\ & + BCE(Classifier(y'), val_X) \end{aligned} \quad (5)$$

390 4.6. Random Grayscaleing

391 In order to force the model to place a greater focus on
 392 geometric features, I added a random grayscale transform
 393 with a decaying random rate,
 394

$$p = 0.25e^{-\frac{curr}{total}} \quad (6)$$

395 Where curr is the current step of the epoch and total is the
 396 total number of steps in the epoch.
 397

406 5. Experiments

407 After initial testing I found that the Class loss and Reten-
 408 tion losses give unremarkable results and actually worsened
 409 output fidelity. I trained the following models demonstra-
 410 ting the efficacy of random grayscaling and the diff loss.
 411 Results will be demonstrated using the following pair of im-
 412 ages from each class: See Figure 1
 413

414 5.1. Baseline

415 The results are shown in Figure 2.
 416

417 5.2. Diff Loss

418 The baseline model was modified to include the Diff
 419 Loss term(3.2) with a lambda weight of 1.
 420

421 The results are shown in Figure 3.
 422

423 5.3. Grayscaleing

424 The baseline model was modified to include the random
 425 Grayscaleing(3.6).
 426

427 The results are shown in Figure 4.
 428

429 5.4. Random Grayscaleing and Diff Loss

430 The baseline model was modified with the addition of
 431 the random Grayscaleing(3.6) and the Diff Loss(3.2).
 432

433 The results are shown in Figure 5.
 434

432
433
434
435
436
437
438
439
440
441
442
443
444457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485Figure 1. Testing Images -
top: Golden Retriever(real dog), bottom:Arcanine(real Pokemon)

6. Risky Experiments

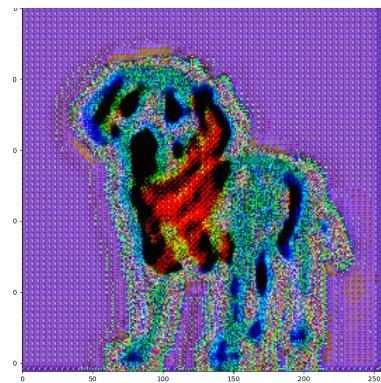
Due to Colab timing out and revoking my GPU access(and also making me lose my trained models!!!), the following models were trained locally and thus were trained for only 4000 steps. My risky experiments centered on altering the CycleGAN architecture as highlighted in 3.5. For this purpose I trained this model with Random Grayscale and the Diff Loss. I also trained a baseline model with the Random Grayscale and Diff Loss for the same number of steps for comparison.

6.1. Random Grayscale and Diff Loss

The baseline model was modified with the addition of the random grayscale(3.6) and the diff loss(3.2). The results are shown in Figure 6.

6.2. Classifier with Random Grayscale and Diff Loss

The classifier model(3.5) was modified with the addition of the random grayscale(3.6) and the diff loss(3.2). The results are shown in Figure 7.

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

7. Summary of Experiments

I performed initial tests on different combinations of my methods and concluded that the Class Loss and Retention Loss were not worth including in further experiments. I then trained 4 models for 10000 steps, focusing on improving the baseline model using Diff Loss and Grayscale. These experiments were greatly successful and performed quite positively. For my risky experimentation, I trained my modified classifier CycleGAN architecture with Grayscale and Diff loss for 4000 steps(due to time and resource constraints) and compared it to the base model with Grayscale and Diff loss trained for the same number of steps. While I am unfortunately unable to see this model perform after extensive training, its early results show pleasantly positive outputs. This is an excellent result for such a risky experiment and is worth further investigation.

8. Analysis

The baseline model(Figure 1) produced extremely unsatisfactory results, performing very little change to the images and producing a heavily blurred quality with a distorted ef-

540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593



Figure 3. Base Model w/ Diff Loss Output(10000) -
top: Golden Retriever(fake Pokemon), bottom:Arcanine(fake dog)

fect on it background. Likely this is due to the model's attempts to replicate the animation style of the anime by placing sharp patterns in the background. The Diff Loss example(Figure 2) made small strides toward geometric feature changes. In the generated dog, it reshapes the mouth of the Pokemon, such that you no longer see the small teeth and makes it seem more like an open mouth of a dog. Its dog to Pokemon conversion adds a lined pattern which masks the realistic fur of the dog.

The Grayscaleing example(Figure 3) produced higher quality outputs which did not necessarily have as drastic changes but seemed less distorted. Notably, it rounded the lines of the fake dog's frame to make it less cartoon like and sharpened the fake Pokemon's frame for the opposite effect. It also attempts to the reduce the prominence of the Pokemon's stripes in it fake dog output.

The combined methods(Figure 4) produced results which retained benefits from both. The Pokemon to dog conversion once again added the open mouth to the fake dog as well as beadier dog eyes. In the dog to Pokemon conversion, the generator added clearly bounded patches of color onto the dog's body, much like a Pokemon's design.

While the combined methods were more effective overall and produced outputs with a balance of quality and drastic changes, the combined method also failed at things its

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



Figure 4. Base Model w/ Grayscaling Output(10000) -
top: Golden Retriever(fake Pokemon), bottom:Arcanine(fake dog)

components succeeded at, such as reducing the prominence on the fake dog' stripes.

The results of these methods on the fake dog image have been compiled in Figure 8 for easy comparison.

Similarly the results of these methods on the fake Pokemon images have been compiled in Figure 9.

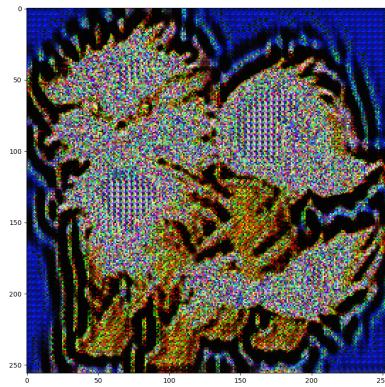
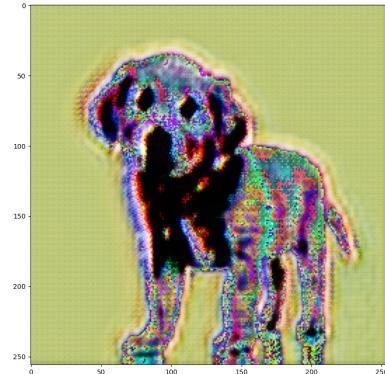
In my further experiments which highlights the classifier model, the output of the classifier model(Figure7) is comparable to the base model(Figure 6) with the same number of steps in quality. This is a pleasant result for a novel architecture as it has not reduced the quality of the model. However, notably, the classifier model after 4000 steps seems to have produced a Pokemon with more defined outline and more solid patches of color than the base model after a similar number of steps. Upon reference to Figure 5, it seems that Figure 7 has identified similar features and made similar changes. Interestingly, it has done this more than Figure 6, which is the same model as Figure 5. While this is an unverified hypothesis, I would speculate that the Classifier model is able to converge more rapidly than the Base model.

9. Conclusion

While we have failed to see significant geometric change using the CycleGAN architecture, we have made very small strides toward improvement. While this could likely be

648
649
650
651
652
653
654
655
656
657
658
659660
661
662
663
664
665
666
667
668
669
670
671672
673 Figure 5. Base Model w/ Grayscale and Diff Loss Out-
674 put(10000) -
675 top: Golden Retriever(fake Pokemon), bottom:Arcanine(fake dog)

676

677 credited to shortcomings in my method and may be possible
678 through some careful tuning of parameters and architecture,
679 the objective of geometric change may be out of reach for
680 CycleGAN. Beyond geometric change alone, this transla-
681 tion task is also quite tricky due to its disagreement with
682 the underlying concept behind unpaired translation. That
683 is, unpaired translation preposes some underlying mapping
684 between the classes. For Pokemon which are only loosely
685 inspired by animals, there likely does not exist such a map-
686 ping between them and dogs. This is unlike the horse to
687 zebra problem where there is a far more realizable strategy
688 for mapping them to each other.689 However, for a difficult unpaired image translation applica-
690 tion, with a very limited dataset, my proposed Diff Loss and
691 Random Grayscale methods produced arguably higher fi-
692 delity outputs than the baseline architecture. Whether or not
693 they have potential for geometric changes, these methods
694 are worth further study for usage in difficult image transla-
695 tion tasks.696 Finally, with very little study, my proposed Classifier Cycle-
697 GAN architecture model shows potential as a more effective
698 variation on the standard CycleGAN architecture.699
700
701702
703
704
705
706
707
708
709
710
711
712
713714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

780

781

782

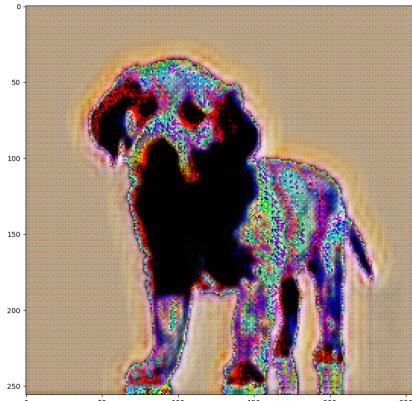


Figure 7. Classifier model w/ Grayscaleing and Diff Loss Output(4000) -
top: Golden Retriever(fake Pokemon), bottom:Arcanine(fake dog)

801

802

803

804

805

806

807

808

809

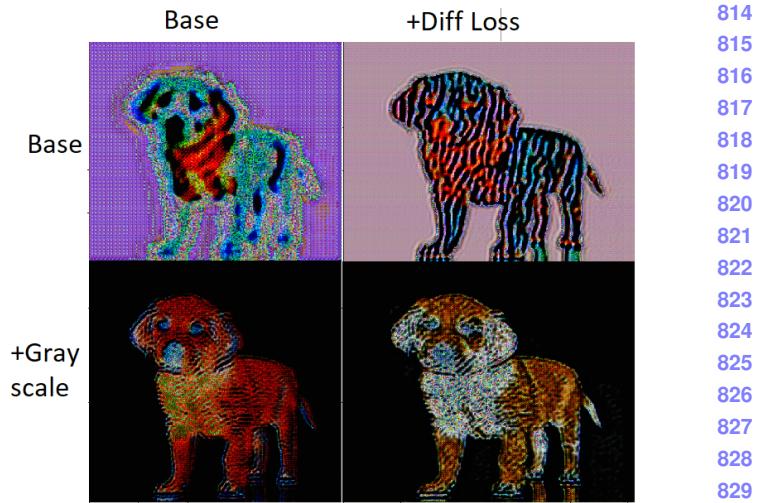


Figure 8. Grid showing results of adding Diff Loss and Grayscaleing to the baseline model on generated pokemon

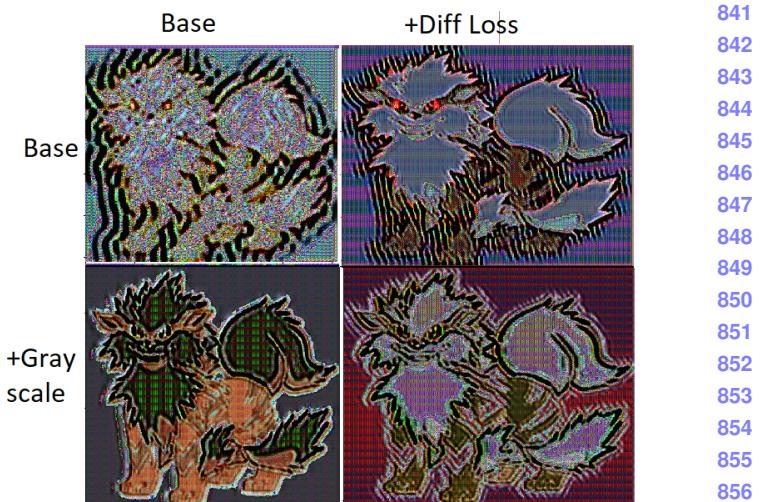


Figure 9. Grid showing results of adding Diff Loss and Grayscaleing to the baseline model on generated dogs