
Exploring: "Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks"

Authors

KTH Royal Institute of Technology KTH Royal Institute of Technology
Altan Senel Ayberk Uckun
altans@kth.se uckun@kth.se

Abstract

1 In recent years, generative adversarial networks have been achieving impressive re-
2 sults, especially in image-to-image translation problems. However, paired datasets
3 are non-existent for many of the tasks, making these networks unattractive to train.
4 CycleGAN successfully addresses this problem by introducing cycle consistency
5 loss. We implement the architecture and the training from scratch and replicate
6 some experiments with the aim of having a deeper understanding of CycleGANs
7 and generative adversarial networks.

8 1 Introduction

9 The image-to-image translation is a problem in computer vision of learning the mapping between
10 an input image domain A and output image domain B [1]. Works on this domain are mostly done
11 by using paired datasets. However, finding paired datasets for two arbitrary image domains is
12 challenging. By paired dataset, we mean that we would need paired images of a cat and a dog that
13 would look exactly similar in their positioning, lighting, background etc.

14

15 This is where CycleGAN, which is a generative adversarial network, excels with learning these
16 mappings from domain A to B with unpaired datasets [1]. Generative adversarial networks are
17 first introduced in 2014, which the network consists of a generative model G and a discriminator
18 model D which estimates whether a sample came from the real training data or was generated by G
19 [2]. Generator and discriminator compete against each other during the learning, hence why these
20 networks are called adversarial. For CycleGAN, there are two generator and discriminator pairs
21 for translations from domain A to B and for the reverse direction. Also in order to ensure proper
22 translation between domains a cycle consistency loss is introduced in CycleGAN for it to allow
23 unpaired image-to-image translation [1].

24

25 Firstly, we tried to understand the architecture in depth and tried implementing CycleGAN from
26 scratch by only using the CycleGAN paper and no supplementary material. However, the paper
27 was not explicit in some parts of the architecture, for example, the spatial padding. After the
28 implementation, we tried to recreate the results for a subset of the experiments in the paper and
29 performed qualitative comparisons with the results from the paper.

30

31 2 Related Work

32 Previously, generative adversarial networks are applied to many challenging vision problems such
33 as realistic image manipulation with user-controlled queries [3], automatic synthesis of realistic
34 images from text [4], predicting future images from a video sequence [5] and many more. Similarly,

35 CycleGAN builds on the concept of adversarial loss introduced in [2]. Adversarial loss in the
36 image-to-image translation context means optimizing the parameters of the generator so that the
37 translated image is indistinguishable from a real image. Discriminators' job is to be able to tell the
38 difference between a translated image and a real image.

39

40 A recent approach in image-to-image translation tackles the pixel-wise image segmentation prediction
41 problem, which can be formulated as an image translation problem, employs fully convolutional
42 networks [6]. CycleGAN builds upon the "pix2pix" framework, which uses a conditional generative
43 adversarial network architecture and a learned loss function instead of a hand-engineered one [7].
44 The main advantage of CycleGAN over the pix2pix framework is that it does not require a paired
45 dataset.

46

47 Previous works on unpaired image-to-image translation include the bayesian approach for inferring
48 the most likely output image [8]. A more similar approach employs a weight-sharing constraint
49 that favours learning a joint distribution of multi-domain images over the product of marginals [9].
50 Different from many other approaches, CycleGAN does not rely on task-specific similarity functions
51 or that input and output should lie in the same low dimensional space [1]. With fewer assumptions,
52 CycleGAN becomes more of a general purpose architecture for different computer vision problems.

53 **3 Data**

54 We used four different datasets and applied the same preprocessing to all of them.

55 **3.1 Datasets**

56 **3.1.1 Horse2Zebra and Apple2Orange**

57 These two datasets consist of images from ImageNet [10]. We downloaded these datasets using the
58 "cycle_gan" collection of TensorFlow Datasets [11] and Tensorflow provides these datasets using
59 the published links for the datasets by the authors of [1]. Authors created these datasets using the
60 keywords "wild horse", "zebra", "apple", and "navel orange" and scaled the images to 256x256
61 resolution. The resulting datasets had 989 horse, 1177 zebra images and 996 apple, 1020 orange
62 images.

63 **3.1.2 Summer2Winter Yosemite**

64 This dataset consists of photos taken in Yosemite National Park and uploaded to Flickr. We accessed
65 the dataset using the UC Berkeley's official directory for CycleGAN Datasets [12]. The authors
66 retrieved the images using the Flickr API with the tag "yosemite" and "data-taken" field, pruned the
67 black and white images and scaled the resolution to 256x256. The resulting dataset had 1273 summer
68 and 854 winter images.

69 **3.1.3 Cityscapes**

70 Cityscapes dataset contains 5000 fine pixel-level annotated and 20k coarse annotated street scenes
71 from 50 different cities taken from recorded video sequences [13]. The authors took 2975 fine
72 annotated images and their raw images as a training dataset and used 500 images of the validation set
73 as the test set. We accessed this dataset using Tensorflow Datasets. Images have 128x128 resolution.

74 **3.2 Preprocessing**

75 Preprocessing includes jittering and normalizing an image but not flipping. To normalize an image,
76 we brought images to the range of [-1, 1] and to jitter, we re-scaled them to 286x286 resolution
77 using the Nearest Neighbor method and randomly cropped the image back to 256x256. This way we
78 augment the images to avoid overfitting.

79 **4 Methods**

80 **4.1 Adversarial Loss**

81 The purpose of unpaired image-to-image translation is to learn a mapping function between two
 82 domains. CycleGAN has two mapping functions $G: X \rightarrow Y$ and $F: Y \rightarrow X$ and also two adversarial
 83 discriminators. G tries to match the generated images' distribution with the target distribution while
 84 D_Y discriminates between real and fake domain images. These two loss components create an
 85 adversarial min-max loss. While the original adversarial loss is consist of the negative log-likelihood,
 86 to better stabilize the training, the authors applied least-square loss following the work [14]. Overall
 87 adversarial loss for the Y domain is given in Equation 1. The first two components of the loss's
 88 right-hand side belong to the discriminator while the third one belongs to the generator. Interchanging
 89 the places of x and y gives other domain's loss.

$$\begin{aligned} \mathcal{L}_{\text{GAN}}(G, D_Y, X, Y) = & \mathbb{E}_{y \sim p_{\text{data}}(y)} [(D_Y(y) - 1)^2] + \mathbb{E}_{x \sim p_{\text{data}}(x)} [D_Y(G(x))^2] \\ & + \mathbb{E}_{x \sim p_{\text{data}}(x)} [(D_Y(G(x)) - 1)^2] \end{aligned} \quad (1)$$

90 **4.2 Cyclic Loss**

91 In theory, any mapping from domain X to Y is correct as long as it matches the target distribution.
 92 Thus, adversarial loss alone does not guarantee a proper image translation. To ensure that only the
 93 required parts of the image change and the image does not lose its authenticity, one must decrease the
 94 number of possible mapping functions by adding constraints. CycleGAN argues that by ensuring
 95 cycle consistency between domains, that is, a translated image by G can be easily translated back
 96 to its original domain by F , we can learn the proper mapping functions. So CycleGAN introduces
 97 Equation 2 as an additional loss to generators.

$$\begin{aligned} \mathcal{L}_{\text{cyc}}(G, F) = & \mathbb{E}_{x \sim p_{\text{data}}(x)} [\|F(G(x)) - x\|_1] \\ & + \mathbb{E}_{y \sim p_{\text{data}}(y)} [\|G(F(y)) - y\|_1] \end{aligned} \quad (2)$$

98 **4.3 Identity Loss**

99 Identity Loss is suggested by [15] and used in CycleGAN for some datasets when required. This
 100 additional loss is added to regularize generators to preserve the original colour of the image. Normally
 101 changing the colour is a valid approach that a generator can apply while, for example, generating
 102 summer-to-winter images. However, in the case of photos-to-paintings, preservation of the colour is
 103 important. This loss is applied in Equation 3 by giving X domain's image to that domain's generator
 104 and expecting the output to be the same.

$$\begin{aligned} \mathcal{L}_{\text{identity}}(G, F) = & \mathbb{E}_{y \sim p_{\text{data}}(y)} [\|G(y) - y\|_1] \\ & + \mathbb{E}_{x \sim p_{\text{data}}(x)} [\|F(x) - x\|_1] \end{aligned} \quad (3)$$

105 **4.4 Instance Normalization**

106 Instance normalization, also known as contrast normalization, is an adaptation of batch normalization
 107 for generative networks proposed by [16]. It can be seen as equal to the batch normalization when the
 108 batch size is 1 but unlike batch normalization, instance normalization is applied also in test time. This
 109 kind of normalization simplifies the generation process by removing the instance-specific contrast
 110 information from the image.

111 **4.5 Reducing Model Oscillation**

112 Oscillation of the loss is a common occurrence while training Generative networks. To make the
 113 discriminator more stable for a better training, [17] suggests using the collection of generated images
 114 instead of using the latest one as a strategy. CycleGAN uses the latest 50 generated images as an
 115 image pool to draw a batch size of images to calculate the discriminator loss. This way discriminator
 116 can continue learning its mistakes from the previous images while adapting to new ones.

117 **4.6 PatchGAN**

118 For the discriminator of the cycleGAN, 70x70 PatchGANs [7] are used. These GANs are designed to
119 classify overlapping 70x70 patches of an image. This way discriminator has fewer parameters, runs
120 faster and can work on arbitrary size images. Decreasing the patch size up to the level of PixelGAN
121 only decreases the quality (makes it blurry) while increasing up to ImageGAN level only makes it
122 harder to train. A PatchGAN discriminator models the image as a Markov random field, assuming
123 independence between pixels separated by more than a patch diameter. [7]

124 **4.7 Upsampling the checkerboard effect**

125 In the generator, after a series of down-sampling operations using 2d convolutions, upsampling is
126 done using deconvolutions to obtain the original spatial dimensions. Our deconvolution kernel is of
127 size 3 and stride is 2, therefore, there is an uneven overlap when sliding these kernels. This effect is
128 called the checkerboard effect or in some cases, mosaic pattern [18]. This effect is especially observed
129 in the early samples generated by the generator. For future work, new forms of regularization or
130 parameter sharing can be employed as suggested in [18]. Alternatively, all deconvolutions can be
131 replaced with a different upsampling technique that is not susceptible to the checkerboard effect.

132 **5 Experiments**

133 **5.1 Horse2Zebra**

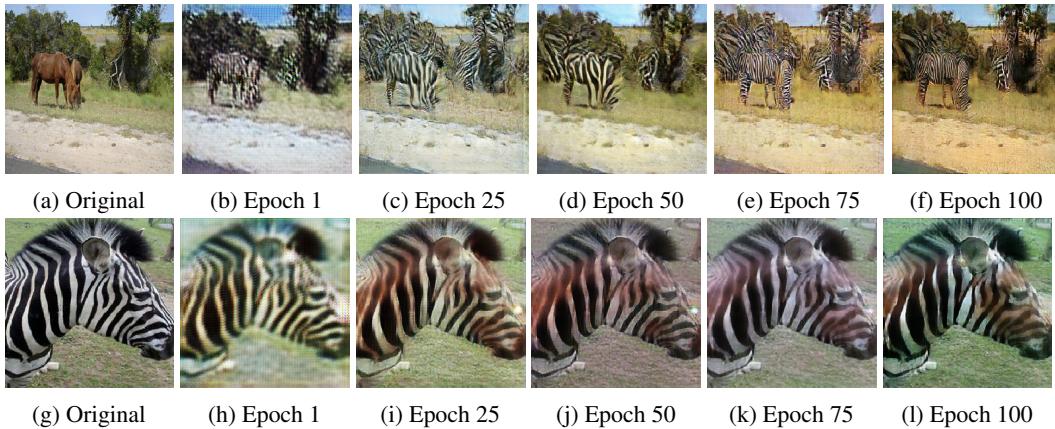


Figure 1: Generated Images using Horse2Zebra Dataset. First row horse-to-zebra, second row zebra-to-horse

134 Figure 1, shows the generated images throughout the training of CycleGAN on Horse2Zebra
135 dataset. We can see that epoch 1 images are all blurry and low quality but it shows a little bit of
136 the characteristics of the target domain, like black and white lines for zebras and brown colour
137 for a horse. On epoch 25, it becomes clearer that the network learned that it should apply these
138 transformations more commonly and also network produces higher quality images. On epoch 75, the
139 network not only learned that it should change the texture and the colour of the animal but also the
140 background since horses generally live around greener grass compared to zebras living in savannas.
141 On epoch 100, it starts solving the problem of mixing the background tree with the horse itself and
142 also produces better quality images.

143

144 We train our models for 100 epochs compared to the original experiments with 200 epochs and
145 also we use mixed precision in calculations because of resource constraints. Our images are not
146 perfect but they show clear signs of learning and somewhat successful image translation. The authors
147 also warn about the existence of failure cases. It is possible to find better-transformed images and
148 worse-transformed images in the dataset thus showing the results turns into a cherry-picking process
149 since there is no proper consistency and a metric to measure the overall performance of the network.

150 We can only visually evaluate these unpaired images and this is what the authors did in the CycleGAN
151 paper using "Mechanical Turks" and by asking real or fake.

152 **5.2 Apple2Orange**



Figure 2: Original apple image



Figure 3: Generated orange image after 20 epochs



Figure 4: Original orange image



Figure 5: Generated apple image after 20 epochs

153 Out of the image-to-image translations tackled in this report, the most trivial mapping seems to be in
154 the case of apple2orange. With only colour and texture changes, it is possible to make an apple look
155 like an orange and vice-versa. This task does not require learning geometric transformations of the
156 objects, therefore with fewer epochs, decent results were achieved by generators.

157 Colour transformation is learned really early on in the training process. With more epochs, generators
158 learn to remove artefacts such as the checkerboard effect and make the images sharper, which will
159 complicate the job of the discriminator.

160 **5.3 Summer2Winter Yosemite**

161 For this dataset, again, we see somewhat more successful results compared to Horse2Zebra dataset
162 since still only a colour change somewhat can solve the problem. Over the course of the training, the
163 network learns that generated summer images must be greener and brighter while winter images
164 must be whiter and darker. The translation from Figure 6a to Figure 6f and from Figure 6m to Figure
165 6r clearly shows this phenomenon.

166

167 Moreover, in Figure 7, we can see a generated winter image using a dark blue and orange colour
168 plate which suggests a more nightly autumn image. This kind of image generation can happen since
169 we are not using identity loss and these kinds of images fit more to the distribution of winter images
170 rather than summer images.

171 Loss graphs of the CycleGAN while training ON Summer2Winter Yosemite are given in Figure 8.
172 Losses are smoothed using 0.99 exponential moving averaging yet the real graph can be seen in the
173 background. The original loss graph has lots of oscillation since the loss values are reported after
174 every 5 batches with batch size 1. Also, the oscillating loss graphs are common in adversarial learning
175 since the generator and the discriminator battle against each other and when one of them is doing
176 better this means the other one is doing worse hence oscillating the loss. We can see that the biggest
177 contributor to the generator loss is the cyclic loss and it decreases steadily. Decreasing its value
178 lowers the total generator loss even though the adversarial loss of the generator slightly increases.
179 Also for the discriminators, we can see that their loss generally stays stable after a quick drop.

180 **5.4 Cityscapes**

181 The original image is the ground truth pixel-wise segmentation of a street image in this experiment.
182 There are two generators that are learning the mapping from pixel-wise segmentation to street
183 photography and reverse mapping. Here we specifically look at the images generated by the generator
184 that creates street photography.

185 As the result of deconvolution with a kernel size of 3 and stride of 2, we can clearly see the
186 checkerboard effect in the generated image after 1 epoch of training. The image is quite blurry
187 and there are not many details. After 50 epochs of training, the generator got better at generating
188 realistic images and the checkerboard effect is minimized. The image is less blurry, however, a good



Figure 6: Generated Images using Summer2Winter Yosemite Dataset. The first two rows are summer-to-winter, third and fourth rows are winter-to-summer



Figure 7: Generated relatively dark blue summer-to-winter image

189 discriminator should still be able to determine if it is generated and not sampled from a real dataset.
 190 The generated image after 100 epochs of training is much more realistic. The road and the trees look
 191 somewhat realistic.

192 Cyclegan is incentivized to create sharp images, as greyish blurred outputs can be identified by the
 193 discriminator network [1]. In this experiment, we clearly see that the generated image got sharper
 194 with more training hence getting closer to the true colour distribution.

195 6 Conclusion

196 We learned that image-to-image translation is an achievable task even without the existence of paired
 197 datasets by learning the domain distributions and proper well-constrained mapping functions between
 198 these distributions. Cyclic loss, identity loss and several other ideas implemented in the CycleGAN is
 199 paper show the way of finding solutions to these kinds of problems. Even though CycleGAN is
 200 successful with translations that include colour and texture changes, it fails with the ones that

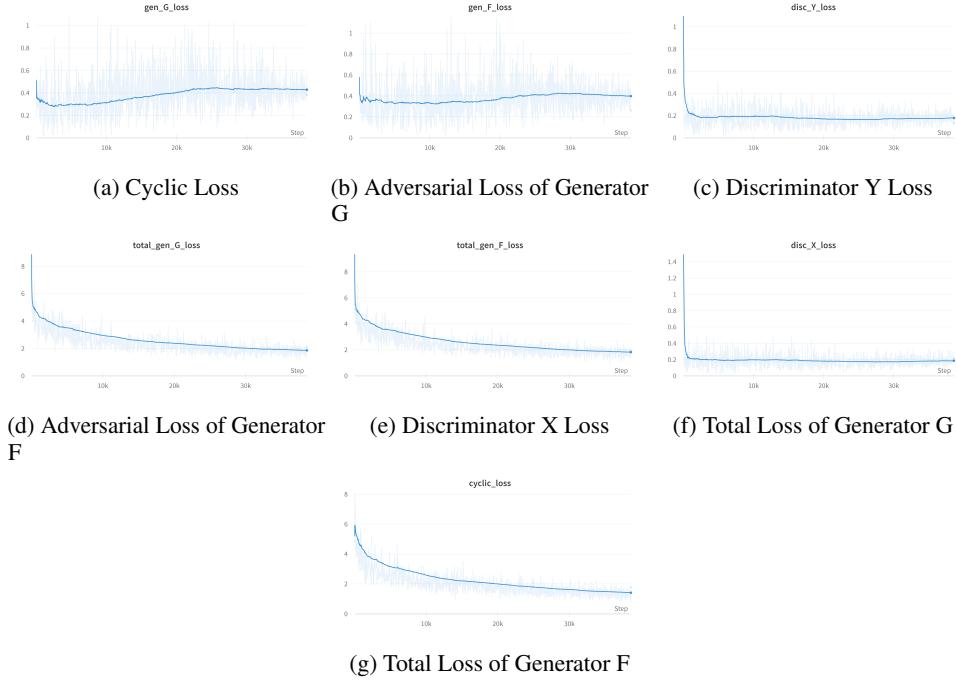


Figure 8: Loss Graphs of CycleGAN Training

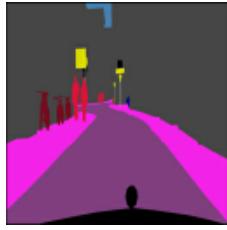


Figure 9: Original segmentation image



Figure 10: Generated city image after 1 epoch



Figure 11: Generated city image after 50 epochs



Figure 12: Generated city image after 100 epochs

201 require geometric changes. This requires further exploration. Moreover, the non-existence of proper
 202 quantitative methods to measure the performance of the network on unpaired domain translation
 203 tasks is a problem which must be considered to create a proper baseline for this task.
 204

205 We also learned that training GANs is a challenging, resource-hungry process. We were able to
 206 reproduce experiments to a reasonable extent, however, this challenge also taught us the ways how
 207 one can efficiently process data and write faster training functions. Implementation from scratch
 208 was a more challenging task than we anticipated thus while we are experimenting with datasets, we
 209 also needed to experiment with the network architecture itself, ie. unspecified ResNet block, left
 210 out activation, ambiguous padding and we needed to experiment with the differentiation library we
 211 use since differentiation of multiple models under the same custom loop, implementing extra image
 212 pooling components and trying to do all these in a faster fashion by applying several tweaks and by
 213 using graph functions for faster execution is not an easy task. Also, an artefact that was discovered
 214 during experimentation was the checkerboard effect, which was the result of our deconvolutions in
 215 the generative model. Instead of the deconvolutions mentioned in the paper, for future work, different
 216 upsampling techniques can be experimented with to get rid of the checkerboard effect.

217 You can find the relevant written-from-scratch code, network architecture and imple-
 218 mentation details in our GitHub repository: [https://github.com/ayberkuckun/
 219 DeepLearningDataScienceProject](https://github.com/ayberkuckun/DeepLearningDataScienceProject)

220 **References**

- 221 [1] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. “Unpaired image-to-image trans-
222 lation using cycle-consistent adversarial networks”. In: *Proceedings of the IEEE international*
223 *conference on computer vision*. 2017, pp. 2223–2232.
- 224 [2] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil
225 Ozair, Aaron Courville, and Yoshua Bengio. “Generative adversarial nets”. In: *Advances in*
226 *neural information processing systems* 27 (2014).
- 227 [3] Jun-Yan Zhu, Philipp Krähenbühl, Eli Shechtman, and Alexei A Efros. “Generative visual
228 manipulation on the natural image manifold”. In: *European conference on computer vision*.
229 Springer. 2016, pp. 597–613.
- 230 [4] Scott Reed, Zeynep Akata, Xinchen Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak
231 Lee. “Generative adversarial text to image synthesis”. In: *International conference on machine*
232 *learning*. PMLR. 2016, pp. 1060–1069.
- 233 [5] Michael Mathieu, Camille Couprie, and Yann LeCun. “Deep multi-scale video prediction
234 beyond mean square error”. In: *arXiv preprint arXiv:1511.05440* (2015).
- 235 [6] Jonathan Long, Evan Shelhamer, and Trevor Darrell. “Fully convolutional networks for
236 semantic segmentation”. In: *Proceedings of the IEEE conference on computer vision and*
237 *pattern recognition*. 2015, pp. 3431–3440.
- 238 [7] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. “Image-to-image translation
239 with conditional adversarial networks”. In: *Proceedings of the IEEE conference on computer*
240 *vision and pattern recognition*. 2017, pp. 1125–1134.
- 241 [8] Rómer Rosales, Kannan Achan, and Brendan J Frey. “Unsupervised image translation.” In:
242 *iccv*. 2003, pp. 472–478.
- 243 [9] Ming-Yu Liu and Oncel Tuzel. “Coupled generative adversarial networks”. In: *Advances in*
244 *neural information processing systems* 29 (2016).
- 245 [10] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. “ImageNet: A large-
246 scale hierarchical image database”. In: *2009 IEEE Conference on Computer Vision and Pattern*
247 *Recognition*. 2009, pp. 248–255. DOI: 10.1109/CVPR.2009.5206848.
- 248 [11] *Cycle_gan: tensorflow datasets*. URL: https://www.tensorflow.org/datasets/catalog/cycle_gan.
- 249 [12] *Index of /taesung_park/CycleGAN/datasets*. URL: https://people.eecs.berkeley.edu/~taesung_park/CycleGAN/datasets/.
- 250 [13] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo
251 Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. “The Cityscapes Dataset for Semantic
252 Urban Scene Understanding”. In: *Proc. of the IEEE Conference on Computer Vision and*
253 *Pattern Recognition (CVPR)*. 2016.
- 254 [14] Xudong Mao, Qing Li, Haoran Xie, Raymond Y.K. Lau, Zhen Wang, and Stephen Paul Smolley.
255 “Least Squares Generative Adversarial Networks”. In: *2017 IEEE International Conference on*
256 *Computer Vision (ICCV)*. 2017, pp. 2813–2821. DOI: 10.1109/ICCV.2017.304.
- 257 [15] Yaniv Taigman, Adam Polyak, and Lior Wolf. *Unsupervised Cross-Domain Image Generation*.
258 2016. DOI: 10.48550/ARXIV.1611.02200. URL: <https://arxiv.org/abs/1611.02200>.
- 259 [16] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. *Instance Normalization: The Missing*
260 *Ingredient for Fast Stylization*. 2016. DOI: 10.48550/ARXIV.1607.08022. URL: <https://arxiv.org/abs/1607.08022>.
- 261 [17] Ashish Shrivastava, Tomas Pfister, Oncel Tuzel, Josh Susskind, Wenda Wang, and Russ Webb.
262 *Learning from Simulated and Unsupervised Images through Adversarial Training*. 2016. DOI:
263 10.48550/ARXIV.1612.07828. URL: <https://arxiv.org/abs/1612.07828>.
- 264 [18] Jon Gauthier. “Conditional generative adversarial nets for convolutional face generation”.
265 In: *Class project for Stanford CS231N: convolutional neural networks for visual recognition*,
266 *Winter semester 2014.5* (2014), p. 2.