

000
001
002
003
004
005
006
007
008
009
010
011054
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

CycleGan: Transfer of style from image to Monet

Meiyi Luan

Xavier Morin

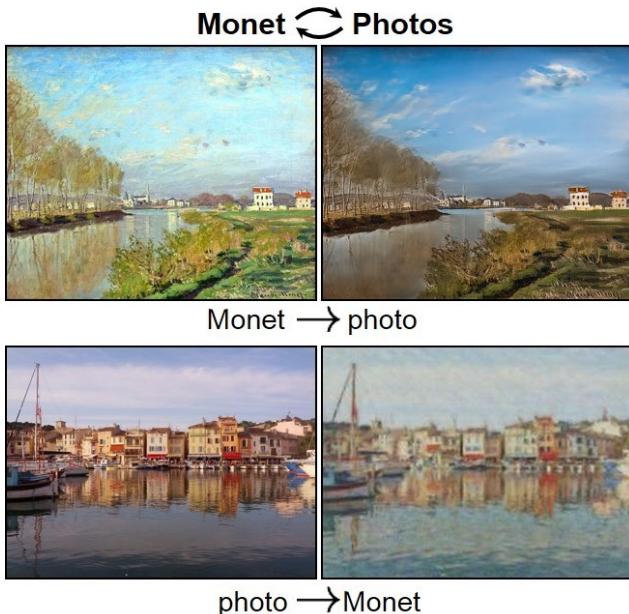


Figure 1: This figure represents our goal [3].

Abstract

We chose our project from Kaggle: "I'm Something of a Painter Myself"[2]. Our goal was to generate paintings in the style of Monet. We had in mind to use a GAN, by training a Generator to map a latent space to an image. This goal has been simplified to transferring the style of Monet to existing images. We used Image-to-Image translation with Cycle-Consistent Generative Adversarial Network[3]. We did not achieve the resolution or quality that we aspired to. However, we learned a lot and are satisfied with the results.

1. Introduction

We were amazed by the various accomplishments of GANs before taking this course, for its applications in Science to improve astronomical images or to generate new molecules, or its use in Art and Fashion. The whole idea

behind GANs is also thrilling, having two networks competing one against other. This said we thought that we could train a GAN to generate paintings of Monet from ground up using Progressive Growing GAN (See our tentative to train a Progressive Growing GAN at the end of this document). It was rather naïve, and we rapidly changed our direction to use a Cycle GAN in order to do image translation like in Figure 1. We give all the credits to **Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks**[3].

2. Theory

Cycle GAN

A GAN is composed of two Neural Networks: a Discriminator and a Generator. Which are trained simultaneously to compete against each other. A Cycle GAN is composed of two GANs : two Discriminators and two Generators. One Gan is learning the mapping $G : X - \rightarrow Y$ and the other the mapping $F : Y - \rightarrow X$. The idea behind using two GANs is to be able to introduce a second parameter to the loss function. Here are the default loss functions for a GAN:

$$L(D) = \frac{1}{2m} \sum_{i=1}^m [(D(x_i) - 1)^2] + \frac{1}{2m} \sum_{i=1}^m [(D(G(z_i)))^2]$$

$$L(G) = \frac{1}{m} \sum_{i=1}^m [(D(G(z_i)) - 1)^2]$$



Figure 2: Translating Monet to Picture without Cycle Loss

Those loss functions basically cover the characteristics of the style, the Discriminator recognizes Monet's style or not and the Generator does well if it can convey Monet's style. However, with only those loss functions, the Genera-

108 tor generally will not preserve the scenery of the input image.
 109 Since we are only asking it to convey the style not to
 110 preserve the original image structure, we need to introduce
 111 the cycle loss.
 112

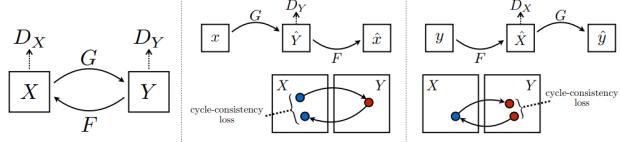


Figure 3: The Cycle Loss [3]

123 The cycle loss will assure that the Generator changes
 124 the style of the image and also preserves the original image
 125 scenery. Finally, we also applied an Identity loss, whose
 126 function will tell the generators to let an image unchanged
 127 if it is already in the desired domain. It also helps to
 128 preserve the colors of the input when applying the full
 129 cycle. Here are those two extra losses:
 130

$$Cycle_L(G) = \frac{1}{m} \sum_{i=1}^m (y_i - G_{x \rightarrow y}(G_{y \rightarrow x}(y_i)))^2$$

$$Identity_L(G) = \frac{1}{m} \sum_{i=1}^m (y_i - G_{x \rightarrow y}(y_i))^2$$

Residual Block

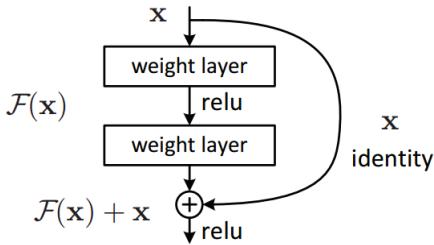


Figure 4: A Residual Block

150 There is a limit to how many layers can be stacked
 151 in a neural network. When the error is propagated from
 152 the last to the first layer, it slowly vanishes. If the neural
 153 network is too big, the first layer would never get updated
 154 and would never learn anything. This is known as the
 155 vanishing gradient problem. The Residual Block proposes
 156 a solution to this problem, by letting the input flow through
 157 the network unchanged. This elegant solution, allows to
 158 train bigger networks and also lets the network directly
 159 learn the Identity function.
 160

ResNet

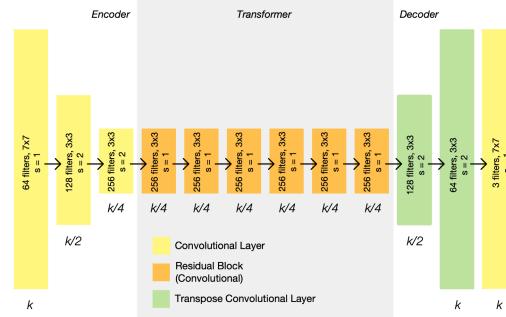


Figure 5: ResNet Generator Architecture .

177 The ResNet Generator applied in [3] makes use of
 178 the Residual Block. It is composed of an Encoder, a
 179 Transformer and a Decoder. The Encoder is composed of
 180 three Convolutional layers that each reduces the size of the
 181 Image by a factor of 4. The transformer of the image is
 182 composed of a series of Residual Blocks, 6 for 128x128
 183 and 9 for 256x256 images. The Decoder brings back the
 184 image to its original form performing deconvolution three
 185 times.
 186

Patch GAN

189 A Patch GAN is a type of Discriminator that only pen-
 190 nalizes structures at the patch level. Each patch is a
 191 NxN subImage taken from the Image. The Discriminator
 192 evaluates if each patch is part of the desired Domain. To
 193 calculate the loss, we simply take the mean of all those
 194 patches. This type of Discriminator has been proven to
 195 be highly effective to discriminate different types of texture.
 196

U-Net

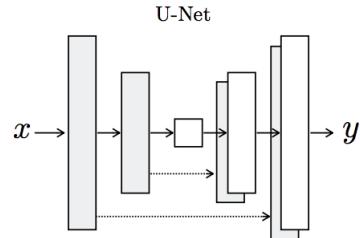


Figure 6: Architecture of the U-Net Generator Model .

209 The U-Net Generator is composed of Encoders and De-
 210 coders that meet at a bottleneck. The Strength of the U-Net
 211

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

is that each Encoder is directly connected to its correspondent Decoder, the first Encoder with the Last Decoder and so on. These Generators rely on Convolution, Deconvolution, LeakyReLU and ReLU activation layers and Batch-Normalization.

3. Methodology & Experimental Results

We wanted to start simple with Cycle GAN and found an assignment from the University of Toronto on the subject, and did it. We learned the basics with this remarkably simple architecture.

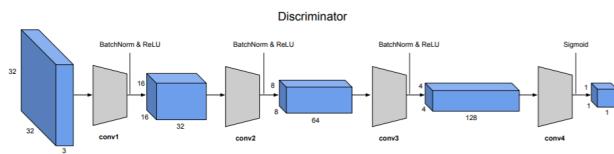


Figure 7: The Discriminator [1].

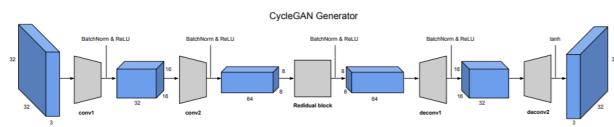


Figure 8: The Generator [1].

Note that this implementation takes input images of size 32x32 and that the discriminator is not a patch GAN, it simply returns 0 or 1. Our first trial was with Emojis (Figure 9). We tried to convert Apple's Emojis to Windows' Emojis.



Figure 9: Translating emojis from Apple to Window's style

After learning the basics, we started working with the dataset provided by Kaggle. During all of our experiments we kept some factors constant, following the recommendations found in [3]. We used the Adam solver, we set lambda=10 and all networks have been trained from scratch

with a learning rate of 0.0002. For most of the experiments we used Batch Normalization with Batch size 16. Have a look at the timeline.



Figure 10: Left: Input X, Center: Ouput G(X), Right: Out-
put F(G(X)). From top to bottom: Scaled Archictecture,
Scaled Architecture and Patch GAN, U-Net, ResNet

In our first attempt, we adjusted the architecture used with Emojis and used 256x256 images as inputs. We kept exactly the same logic but with seven convolutional layers in the Discriminator to reduce the input image to an 1x1 output. We also added two extra layers to the Generator, one convolutional and one deconvolutional. The architecture of the Discriminator made it almost impossible for it to differentiate real images from the fake. The translations of the Generator were not satisfying, mainly because the Discriminator was not giving correct directions. However the Generator was able to perform the cycle and the identity transformations.

That is when we discovered the Patch GAN and decided to introduce it. We implemented a 70x70 Patch GAN with a total of five layers of convolutions, resulting in a 30x30 out-

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

324 put. But introducing the new Discriminator wasn't enough
 325 and it even seemed to decrease the quality of the result. We
 326 concluded that the Discriminator was too powerful com-
 327 pared to the Generator.
 328

329 When searching for a new architecture we encountered
 330 U-Net Generator model, which finally gave us our first con-
 331 vincing translation. We did a lot of experiments with the
 332 U-net Generator. We met lots of funny results along the
 333 way. I invite you to look at the end of the document for our
 334 translation Monet - \downarrow Monet on Fire. Wait, wait, finish to
 335 read our report before going at the end ...

336 Finally, we wanted to try the architecture of the Gener-
 337 ator used in the original paper. Our best results came from
 338 the ResNet. We changed to use Instance Normalization with
 339 Batch size 1. With this architecture we also introduced the
 340 concept of a Pool for fake Images. This way the Discrimi-
 341 nator was trained on both past and new Generated images,
 342 making the learning more stable. Note that the last image
 343 is not good, because the Generator from Monet to picture
 344 broke when we saved it.
 345

346 We kept this architecture until the end. To improve our
 347 result we did some data augmentation. We wanted to three-
 348 fold the number of Paintings of Monet, so we flipped the
 349 Paintings horizontally and vertically.
 350
 351

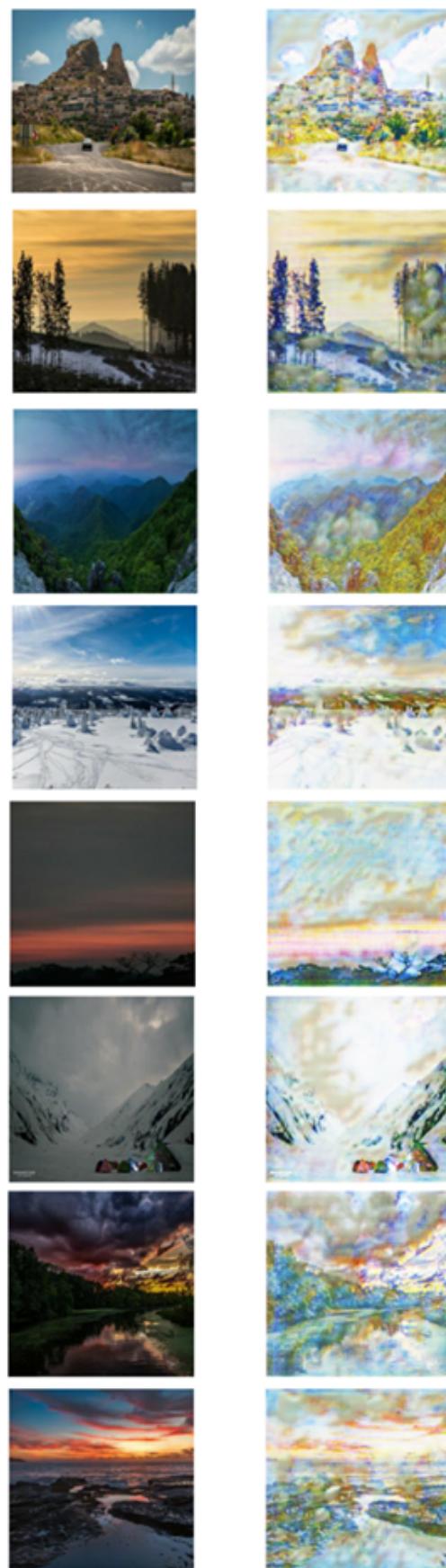


352
 353
 354
 355
 356
 357
 358
 359
 360 Figure 11: Data Augmentation
 361
 362

363 Our last endeavour was with scheduling. We trained the
 364 Models over 70 000 images with a constant learning rate
 365 and again over 70 000 images with a linearly vanishing
 366 learning rate to zero.
 367
 368
 369

4. Conclusion

370 We are very glad to have chosen this project, we didn't
 371 really think that it was possible. We learned enormously
 372 about PyTorch and how to conduct a Neural Network
 373 project. We are overall satisfied with our submission, but
 374 we plan on training the network again in the Winter vaca-
 375 tion and submit our project to Kaggle.
 376
 377

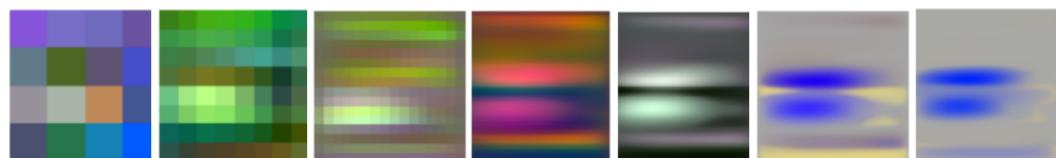


378
 379
 380
 381
 382
 383
 384
 385
 386
 387
 388
 389
 390
 391
 392
 393
 394
 395
 396
 397
 398
 399
 400
 401
 402
 403
 404
 405
 406
 407
 408
 409
 410
 411
 412
 413
 414
 415
 416
 417
 418
 419
 420
 421
 422
 423
 424
 425
 426
 427
 428
 429
 430
 431

432 **References** 486
433
434 [1] Cyclegan a4 from university of toronto. 3 487
435 [2] I'm something of a painter myself. 1 488
436 [3] Jun-Yan Zhu Taesung Park Phillip Isola Alexei A. Efros. Un- 489
437 paired image-to-image translation using cycle-consistent ad- 490
438 versarial networks. 2017. 1, 2, 3 491
439
440 **Appendix: Extra Results (Optional)** 492
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457



458
459 Figure 13: Monet to Monet on Fire 513
460
461
462
463



464
465
466
467
468
469
470
471
472 Figure 14: Progessive Growing GAN, from 4x4 to 256x256 526
473
474
475
476
477
478
479
480
481
482
483
484
485