

Summer Research Program in Industrial and Applied Mathematics



SEOUL
NATIONAL
UNIVERSITY

Sponsor
⟨Tencent⟩

Final Report

⟨Sketch to Image Generation⟩

Student Members

⟨Seon Gyu PARK⟩ (Project Manager), ⟨SNU⟩,
⟨SunQ0313@gmail.com⟩
⟨Chun Wai WONG⟩, ⟨HKUST⟩
⟨Fuzhe LI⟩, ⟨HKUST⟩

Academic Mentor

⟨Dr. Ningchen YING⟩, ⟨mancying@ust.hk⟩

Sponsoring Mentors

⟨Prof. Yu-Wing TAI⟩, ⟨yuwingtai@tencent.com⟩

Consultants

⟨Name⟩
⟨Name⟩

⟨Date: 7th August 2018⟩

Abstract

In this project, we investigated sketch-to-image translation by implementing CycleGAN to learn the mapping between human face sketch to realistic photograph. We used U-net to form the Generative Adversarial Networks This makes the model possible to train end to end from very few images and guaranteeing the performance of the mapping.

Acknowledgments

It is appropriate in the Acknowledgments to thank individuals or organizations who made especially noteworthy contributions to your project. Elsewhere, within the body of the report, you can acknowledge more specific contributions where appropriate. These are matters of courtesy and professional ethics. As an example:

The RIPS L^AT_EX report template has been developed by Mike Raugh with advice and assistance from Oleg Alexandrov and Shawn Cokus in the early stage of development and general support of IPAM and the System Administration staff. The first RIPS template was based on an early version of the Math Clinic's report template at Harvey Mudd College; there the original template has been improved and is managed by Claire Connelly, the HMC Math Department's system administrator. Claire and her co-authors offer coding advice, a wealth of references, and a note about the origin of the template in their current edition, the `sample-clinic-report.pdf` accessible at <http://www.math.hmc.edu/computing/support/tex/sample-report>. Claire copyedited the third edition of Grätzer's *Math into L^AT_EX*, most of which work seems to have survived into the fourth edition: *More Math into L^AT_EX* [?].

When acknowledging individuals in this section, it is OK to use the names by which you know and speak to them. Here it is OK to write "Oleg Alexandrov." But you must be formal on the Title page and elsewhere within the report, where it is proper to specify honorifics, e.g., Dr. or Prof. On the Title page you would write "Dr. Oleg Alexandrov," and likewise within the body of the report if you were acknowledging him for a specific contribution, Claire Connelly uses no honorific, so you would use just her name on the title page. When in doubt, check the person's business card or follow usage on the person's web page.

As a result of suggestions from users, this Sample Report and its source are under continual improvement. Please contact the RIPS program director for your suggestions. An up-to-date list of changes is recorded in the "Revisions" folder for the Master Template Folder.

Contents

Abstract	3
Acknowledgments	5
1 Introduction	13
1.1 Project Goal	13
1.2 Related Works	13
2 Background	15
2.1 Related Models	15
3 Processing Data	17
3.1 Collection and Refinement of Data	17
4 Model Structure	21
4.1 Network Architectures	21
4.2 Stabilization of GAN Training	21
5 Result	23
6 Conclusion	27
6.1 Acknowledgements	27
7 Reference	29
8 Appendix	31
APPENDIXES	
A BibT_EX Sample Records, Record Types and Fields	33
B Where to find this sample RIPS report?	35
C Glossary	37
D Abbreviations	39
REFERENCES	

List of Figures

3.1 Preprocessing examples. Original images containing faces(top) and results after preprocessing steps(bottom)	18
5.1 Original input images(left) and result of translations(right) by networks trained before input data is not aligned. Image size is 256 by 256 and model is trained for 64 epochs.	23
5.2 Original input images(left) and result of translations(right) by networks trained on datasets after face alignment, but before number of smiling faces was reduced. Image size is 128 by 128 and model is trained for 128 epochs.	24
5.3 Original input images(left) and result of translations(right) by networks trained on datasets after face alignment and number of smiling faces was reduced. Image size is 128 by 128 and model is trained for 128 epochs.	25

List of Tables

3.1 Number of images used for training in each component.	17
---	----

Chapter 1

Introduction

1.1 Project Goal

The goal of this project is to build a system that can generate photo-realistic images from rough sketch pictures. To serve this purpose, we utilized the Cycle-GAN [13] with sketch images collected from google image search and Celebrity A dataset [7]. The detailed information about network architecture and processing the collected data will be provided at 4 and at 3 respectively.

1.2 Related Works

Sketch-to-Photograph generation is a sub-problem of Image-to-Image translation where the goal is to learn the mapping between distinct domains of image.

In 2017, Jun-Yan Zhu et al. achieved transforming a horse image into a zebra using Cycle-Consistent Adversarial Networks. However, in the Sketch-to-Image generation, the mapping between photo and sketch always cause loss in dimension of data, resulting in losing the uniqueness of mapping. Our goal is to train a Cycle-Consistent Adversarial Networks to learn a mapping $G : X \rightarrow Y$ such that the distribution of images from $G(X)$ indistinguishable from the distribution Y . In other word, it can generate high quality human face photograph from face sketch image.

GAN: A method of handling Image translation problem is to adopt a Generative adversarial network. Implementing adversarial loss, GAN is able to train the convergence algorithm so that the generator's distribution converges to the data[3]. This contributed to the success of GAN in generating realistic images with an impressive result.

WGAN: As an improvement of GAN, it replaces Jensen-Shannon divergence by Wasserstein distance to calculate the distance between distributions. It solves the mode collapse problem and it proposes an index standing for the quality of training where a larger index represents better results.

Pix2Pix: The framework of Pix2Pix is based on GAN. Adopted paired training examples, Pix2Pix is able to deal with Image-to-Image translation that has similar performance of CycleGAN.

CycleGAN: As an improvement of Pix2Pix, the model of CycleGAN does not

required paired training examples to train the model. It contains two mapping functions $G : X \rightarrow Y$ and $F : Y \rightarrow X$ from two data sets X, Y , and associated with two discriminators D_X and D_Y . In the model, two cycle consistency losses (Forward and Backward) was introduced so that it can achieve two mapping where $F(G(x)) \approx x$ and $G(F(y)) \approx y$ [13]. In this work, we implemented the CycleGAN to work on the Sketch-to-Image translation. Using the U-net instead of the ResNet, we trained the model to specialized in generating realistic human face photograph from a human face sketch as well as generating sketch from human face photograph such that both sketches and human face photograph are indistinguishable from target domain.

LSGAN: Use least square loss function to replace the loss function of GAN, which improved the stability of GAN training and the quality of generated images.

The team would like to thank *Tencent* for the generous sponsorship of this project. Our code is available at https://github.com/SunQpark/SPIA2018_cycle_GAN.

Chapter 2

Background

2.1 Related Models

There are two main approaches for image generation tasks: generative adversarial nets and variational autoencoder. In this project, only GAN-related methods will be considered and implemented. The original GAN proposed by Goodfellow can train a generator to learn the real image distribution by training a discriminator together, which can tell the difference between real images and fake images.

In the field of image translation trained with unpaired image datasets, 3 models, CycleGAN, DiscoGAN, DualGAN are doing the same thing from high level perspective. They train two generative models to learn the mapping between two image domains by training two discriminators at each domain together and considering the cycle-consistency loss, which ensures ignorable difference between images and reconstructed images.

In detail, the differences are CycleGAN uses instance normalization, patchGAN discriminator. To stabilize the training, use least square GAN. Replay buffer no random input z no drop out; L1 distance as cycle consistency

DualGAN use generator and discriminator of pix2pix; no random input z but implemented drop out; use wgan; L1 distance as cycle consistency

DiscoGAN generator: conv, deconv and leaky relu; discriminator:conv+leaky relu L2 distance as cycle consistency

Chapter 3

Processing Data

3.1 Collection and Refinement of Data

In this section, we will explain about the progresses of collecting and refining data we've gone through. Table 3.1 shows the composition of dataset we used in experiments. Unless mentioned otherwise, every dataset used in experiments composed as the proportion as Set1 in table.

Dataset	Sketch			Photo
	Collected	CUFS	Sketch Filter	Celeb A
Set1	920 (30.8%)	571 (19.2%)	0 (0.0%)	1500 (50.1%)
Set2	920 (16.8%)	571 (10.4%)	1000 (18.2%)	3000 (54.6%)

Table 3.1: Number of images used for training in each component.

More details on each sources of data will be provided in following sections.

3.1.1 Photograph

For realistic photographs of human faces, we used Celebrity A [7] which contains more than 200,000 images of celebrity in the world. In addition to having clean images of human faces, this dataset contained ‘attributes’, some additional imformation of images. Though those informations were not used directly as input to our neural networks, they were useful in some progress of preprocessing. We will describe about the details in 3.1.3.

3.1.2 Sketch

The main portion of sketch image dataset we used was collected from Google Image search engine. However, the collected images were highly inconsistent, which was undesirable. figure shows some bad examples of image which were included in the first collection. Thus, these images has gone through preprocessing steps, which will be described in next section 3.1.3 After preprocessing steps, these images were aligned as fig 3.1

We also used CUHK(The Chinese University of Hong Kong) face sketch database(CUFS) [12], containing about 500 face sketches. Although this dataset also provides face photographs paired with the sketch images, we only used the sketch images to train our model to keep it from overfitting to these paired images.

Examples of images in each dataset can be seen in fig.

3.1.3 Preprocess

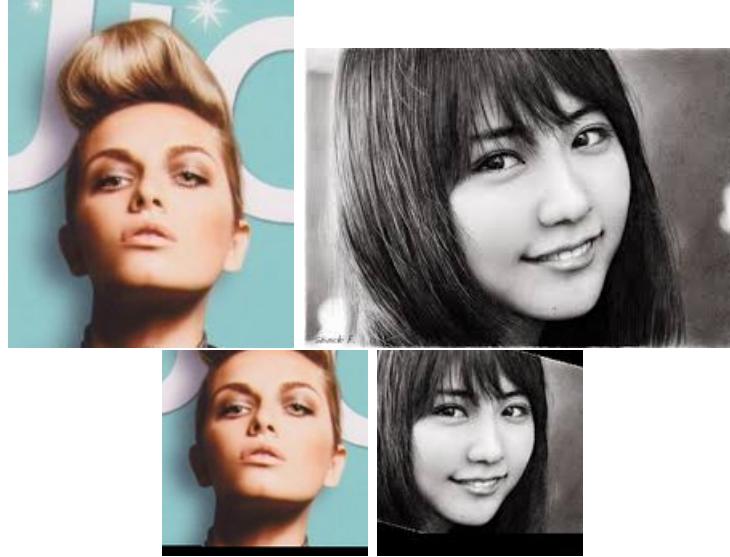


Figure 3.1: Preprocessing examples. Original images containing faces (top) and results after preprocessing steps (bottom)

In this section we deals with the preprocessing steps on face sketches and photos, which were applied separately before the training of model. By experiments of a few weeks, we found this part crucial to the performance of model. The first step of data refinement was to detect every faces for each image. By excluding images not containing recognized faces, we could get rid of 20% of bad examples in sketch database. Then, we applied ‘facial landmark detection’ for the detected faces, in order to get further information concerning the direction of head in the images. We utilized the ‘dlib’ framework [6] for those two steps. Next, we rotated the images to have faces aligned vertically using the detected locations of two eyes. Finally, the images were cropped to have faces in center and 30% of the width of head margins on both sides. The cropped images have 128 by 128 size in the end of preprocessings. Those steps were applied to both sketch and photo images to make them consistent in any attributes other than it is sketch or photo.

After some experiments we found that the model was learning to put smile on the generated photo from sketch images even when the original sketch images were not smiling. This seemed to be originated from the fact that images in photograph have smile in most case, while sketch images haven’t. Since this was undisired effect, we reduced the ratio of smiling faces in the photograph dataset. This process was done

by checking attributes file which Celebrity A dataset provided to see whether given image is smiling or not and using only one images per 20 smiling images.

3.1.4 Generating Sketch Images

Chapter 4

Model Structure

4.1 Network Architectures

As we mentioned earlier, baseline structure mainly used in this project is Cycle-GAN. The Cycle-GAN consists of two GAN models, which learn to translate images from one domain to the other, one direction per each respectively.

Each generator consists of convolutional layers as decoder layers and deconvolutional layers as encoder layers and bottleneck layers. Pooling layers and batch normalization are also implemented in decoders and encoders. At the beginning, we implemented resnet to connect the decoder and encoder as suggested by the original paper. Later, we configured the generators to be u-shape nets where skip connections between downsampling layers and upsampling layers are enabled. In this way, low-level information is shared with upsampling layers without passing through bottleneck layers, which significantly reduces infomation loss and ensures feedback on internal structures of both inputs and outputs.

”For the discriminator networks we use 70 70 PatchGANs , 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 be applied to arbitrarily-sized images in a fully convolutional fashion .”

4.2 Stabilization of GAN Training

wGAN

dualGAN

Instance / Batch normalization

Chapter 5

Result

Figure 5.1 shows some examples of result before the data refinement steps were applied. The results seems to contain some face-like objects, but the original faces in the sketch images were not properly translated, resulting in weird colorization in the result. The result of photograph to sketch translation also shows severe artifacts due to the mismatch of locations of faces between datasets.



Figure 5.1: Original input images(left) and result of translations(right) by networks trained before input data is not aligned. Image size is 256 by 256 and model is trained for 64 epochs.

The result seems to be improved considerably with preprocessing steps. Fig 5.2

shows aligned faces helped the model tell the face and background apart, resulting it to generate better volumes, colors on the paintings.



Figure 5.2: Original input images(left) and result of translations(right) by networks trained on datasets after face alignment, but before number of smiling faces was reduced. Image size is 128 by 128 and model is trained for 128 epochs.

However, the results still looks not realistic, mostly due to mistakenly generated smile on the translated images. After some investigation we could find out that is due to the difference of distribution of train data that people smiled much more on the photos rather than on the paintings. After reducing the ratio of smile in the photograph dataset used as in previous chapter 3, we could get following results 5.3.

While some of the sketch-to-photo results seems quite realistic(3rd image of 2nd column, 2nd, 3rd of 3rd), images generated from low-quality inputs does not seems to be properly translated(2nd image of 2nd column, last ones in 3rd, 4th column). On the other hand, despite the general quality looks fine, some of results(1st, 2nd diagonal positions) resembles average sample images of CUFS datasets rather than contents of input.

Considering the quantity and quality of sketch images collected from google image search, this amount of dependency on input might be considered acceptable.

We have also tried to replace batch normalization by instance normalization to improve the training quality, but the results are not so satisfactory as shown in The mode collapse problem indicates that instance normalization causes instability of training in this case.

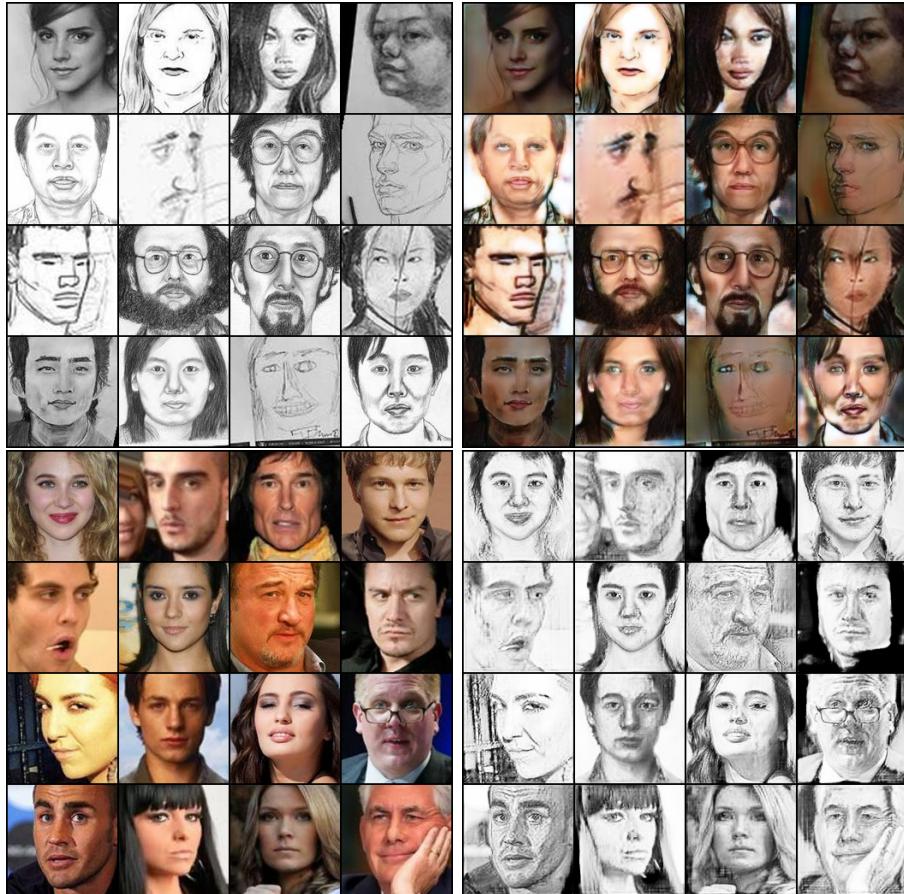


Figure 5.3: Original input images(left) and result of translations(right) by networks trained on datasets after face alignment and number of smiling faces was reduced. Image size is 128 by 128 and model is trained for 128 epochs.

Chapter 6

Conclusion

The point was made in the Acknowledgments section of this Sample Report that it is important to credit others whose work you use—it is a matter of professional ethics and courtesy. In addition to acknowledgment of broad assistance or contributions that you put into the Acknowledgments, you may also need to reference more specific contributions elsewhere in your text. Wherever a distinction is needed, make it clear which part of your work you have borrowed or adapted from others, and provide a reference to the source.

6.1 Acknowledgements

Our team would like to acknowledge Professor Yu-Wing TAI and Dr. Ningchen YING for helpful discussion. Professor Yu-Wing TAI offered a lots of helpful suggestions to model development such as implementing U-net as well as suggestion to data collection. Dr Ningchen YING offered a general introduction and outlining to the project and his patience to guide the project. We would also like to thank Professor Shingyu Leung, Professor Avery CHING, HKUST and SNU MATH department for providing computational resources and a sight seeing trip to Macau for stimulating our creativity.

Chapter 7

Reference

Chapter 8

Appendix

Appendix A

BibT_EX Sample Records, Record Types and Fields

Appendix B

Where to find this sample RIPS report?

Read-only L^AT_EX source code for the RIPS Report Template, sample BEAMER slide presentations, and other L^AT_EX supporting materials are available at,

Computer -> IPAM RIPS FOLDER -> on the R Drive under under "Templates-etc"

Your report will be “copyedited”, i.e., edited for conformance to the RIPS *House Style*. For reference, a table of proofreader’s marks that may be used for markup of your draft is included. It was copied from *The Chicago Manual of Style, 16th ed.* (See original source at: www.chicagomanualofstyle.org/tools_proof.html.)

Appendix C

Glossary

Page vs Leaf: In bookbinding, a trimmed sheet of paper bound in a book; each side of a leaf is a **page**.

Opening: The two pages you see when you open a book. The right-hand **page** is the **recto**—and the left-hand page is the **verso**.

Recto: The front side of a **leaf**; in a book or journal, a right-hand page. To **start recto** is to begin on a recto page, as any major section—e.g., title page, table of contents, preface, chapter, appendix, bibliography—normally does. Contrast **verso**.

Verso: The back side of a **leaf**; the **page** on the left-hand side of an **opening**.

Front matter: As applied to this report, the material that appears in the front of the document, including title page, the abstract, acknowledgments, table of contents, list of figures, list of tables, usually numbered with lowercase roman numerals. RIPS reports initiate pagination with 1 in the front matter and proceed throughout with arabic numerals. This variation of usage is allowed because modern typesetting permits easy re-pagination after pages have been added to the front matter, something not easily done—after completion of the main matter—when typesetting was done by hand.

Main matter: The main part of the document, including the appendixes. **Page** numbers start from 1 using arabic numerals if front matter is enumerated using roman numerals.

Back matter: Material that appears at the back of the document, which in our report includes only the Bibliography.

Appendix D

Abbreviations

IPAM. Institute for Pure and Applied Mathematics. An institute of the National Science Foundation, located at UCLA.

RIPS. Research in Industrial Projects for Students. A regular summer program at IPAM, in which teams of undergraduate (or fresh graduate) students participate in sponsored team research projects.

UCLA. The University of California at Los Angeles.

Selected Bibliography Including Cited Works

- [1] S. BARRATT AND R. SHARMA, *A note on the inception score*, 2018.
- [2] I. GOODFELLOW, Y. BENGIO, AND A. COURVILLE, *Deep Learning*, MIT Press, 2016.
- [3] I. J. GOODFELLOW, J. POUGET-ABADIE, M. MIRZA, B. XU, D. WARDE-FARLEY, S. OZAIR, A. COURVILLE, AND Y. BENGIO, *Generative Adversarial Networks*, ArXiv e-prints, (2014).
- [4] I. J. GOODFELLOW, J. POUGET-ABADIE, M. MIRZA, B. XU, D. WARDE-FARLEY, S. OZAIR, A. COURVILLE, AND Y. BENGIO, *Generative adversarial networks*, 2014.
- [5] K. HE, X. ZHANG, S. REN, AND J. SUN, *Deep residual learning for image recognition*, 2015.
- [6] D. E. KING, *Dlib-ml: A machine learning toolkit*, Journal of Machine Learning Research, 10 (2009), pp. 1755–1758.
- [7] Z. LIU, P. LUO, X. WANG, AND X. TANG, *Deep learning face attributes in the wild*, in Proceedings of International Conference on Computer Vision (ICCV), 2015.
- [8] X. MAO, Q. LI, H. XIE, R. Y. K. LAU, Z. WANG, AND S. P. SMOLLEY, *Least squares generative adversarial networks*, 2016.
- [9] O. RONNEBERGER, P. FISCHER, AND T. BROX, *U-net: Convolutional networks for biomedical image segmentation*, 2015.
- [10] T. SALIMANS, I. GOODFELLOW, W. ZAREMBA, V. CHEUNG, A. RADFORD, AND X. CHEN, *Improved techniques for training gans*, 2016.
- [11] D. ULYANOV, A. VEDALDI, AND V. LEMPITSKY, *Instance normalization: The missing ingredient for fast stylization*, 2016.
- [12] X. WANG AND X. TANG, *Face photo-sketch synthesis and recognition*, IEEE Transactions on Pattern Analysis and Machine Intelligence, 31 (2009), pp. 1955–1967.

- [13] J.-Y. ZHU, T. PARK, P. ISOLA, AND A. A. EFROS, *Unpaired image-to-image translation using cycle-consistent adversarial networks*, 2017.