
Image Synthesis Using a Single Robustness Classifier

Shining Yang, Abhijith Tammangari, Minjung Lee
shining.yang@duke.edu, at396@duke.edu, minjung.lee@duke.edu
Duke University

Abstract

1 In this project, we performed several image synthesis tasks (image generation,
2 image inpainting, super-resolution, image-to-image translation, and sketch-to-
3 image translation) using a pre-trained adversarially robust classifier released by
4 Santurkar et al (2019). Due to the limitation in computation resources, we used the
5 ImageNette dataset, which is a subset of 10 easily classified classes from Imagenet.
6 In line with the paper, we found that the features learned by a basic classifier are
7 sufficient for all these tasks, provided this classifier is adversarially robust. The
8 idea for these augmentations is to take the original adversarial model and use L2
9 FGSM attacks in various means to achieve the synthesis tasks.

10

1 Introduction

11 The idea of utilizing adversarial attacks of L2 FGSM on an adversarially trained robust classifier
12 to perform image synthesis tasks was initially proposed by Santurkar et al (2019). The main goal
13 behind this paper is to utilize L2 FGSM with various loss functions and minimize the worst-case loss
14 of the set of perturbed images in order to achieve tasks such as image generation, super-resolution,
15 image in-painting, and image-to-image translation.

16 Based on their findings a simple FGSM attack with a varying number of iterations, loss functions,
17 and alpha values while minimizing the loss for the batch of samples will easily allow us to perform
18 the various synthesis tasks at hand.

19 On top of proposing this idea, Santurkar and his team also provided a GitHub repository containing
20 links in their paper which contain the suggested L2 FGSM attack attributes as well as model weights
21 for the adversarially trained res-net 50 model. Our goal is to initialize a RES-NET50 model and
22 utilize their adversarially trained RESNET50 model weights to load the pre-trained model. Then we
23 will use our implemented iterations of L2 FGSM attacks with the respective loss functions for each
24 task with the goal of minimizing that total loss. Doing so should enable us to achieve the goal of
25 performing all of our image synthesis tasks.

26 Furthermore, in our use-case we utilized the Imagenette dataset which is a subset of the Imagenet
27 dataset. The original dataset that Santurkar's team utilized to train their adversarial model was the
28 Imagenet data set, so the Imagenette dataset will still work for our goals. The Imagenette dataset has
29 a total of 10 classes which were relabeled to match the original Imagenet data classes.

30

2 Image Generation

31 Realistic image generation has long been a challenge in the computer vision field. According to
32 Santurkar et al (2019), an adversarially robust classifier suffices for image synthesis tasks including
33 realistic image generation. The generation process is based on maximizing the class score of the
34 desired class using the robust model. The goal of this maximization is to add relevant and significantly
35 meaningful features of the desired class to the input image. Here in terms of the input image, we

36 utilize a random seed at the beginning in order to generate a diverse set of input samples. In line with
 37 Santurkar et al (2019), we choose a simple but sufficient multivariate normal distribution to fit the
 38 class-conditional distribution. We visualized the conditional seeds for each class in Figure 2.



Figure 1: Seeds Visualization

39 To maximize the class score for our desired classes, we minimize the loss \mathcal{L} of target label y :

$$x = \arg \min_{\|x' - x_0\|_2 \leq \varepsilon} \mathcal{L}(x', y), \quad x_0 \sim \mathcal{G}_y \quad (1)$$

40 where \mathcal{G}_y refers to the class-conditional seed distribution. We adopted the pre-trained model weights
 41 and hyperparameters released by Santurkar et al (2019) and utilized Projected Gradient Descend
 42 (PGD) with an L2 norm constraint as the strategy to minimize loss. This approach enables us to
 43 achieve image generation from conditional-sampled seeds given any desired class. Images generated
 44 by our method at a 224×224 resolution are shown in Figure 2 (see appendix for images on all
 45 10 classes). The generated images are diverse, realistic, and show clear visual patterns toward the
 46 targeted class.

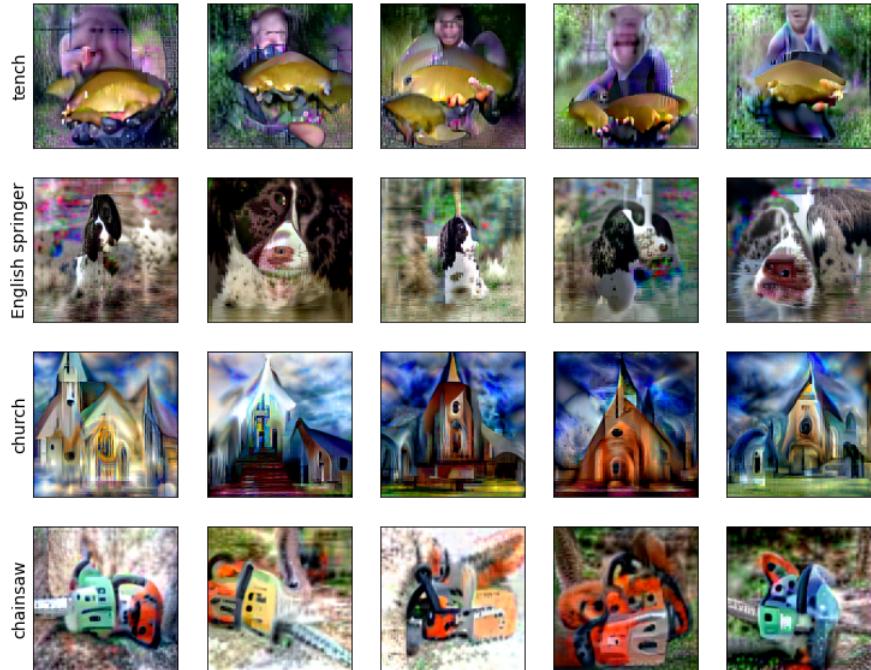


Figure 2: Random samples produced using an adversarially trained robust classifier

47 For comparison, we also implemented a standard non-robust classifier (in this case, a standard
 48 ResNet50) to perform the same tasks. Again, we use the same conditional sampled seeds as input,
 49 minimize the loss using PGD with L2 norm constraint, and visualize the results in Figure 3. The
 50 generated images using a non-robust classifier look vague and noisy without showing any plausible
 51 visual patterns toward the targeted class. The contrast results demonstrate the powerfulness and
 52 sufficiency of a single robust classifier in realistic image generation tasks.



Figure 3: Random samples produced using a standard non-robust classifier

53 3 Image Inpainting

54 The goal of image inpainting is to recover the missing pixels using robust models in a manner that is
 55 perceptually plausible with respect to the rest of the image. Given a robust classifier pre-trained on
 56 uncorrupted data, and a corrupted image x with label y , we solve

$$x_I = \arg \min_{x'} \mathcal{L}(x', y) + \lambda \|(x - x') \odot (1 - m)\|_2 \quad (2)$$

57 where \mathcal{L} is the cross-entropy loss, \odot denotes element-wise multiplication, m is a binary mask, and λ
 58 is a constant (which is 10 here).

59 We used the L2 FGSM attack on the ImageNette dataset, and for an optimization variable x' , we
 60 picked random seeds as 1) pure random noise and 2) sample from class-conditional seed distribution,
 61 the same as used in Image Generation.

62 **Pure random noise** At first, we used pure random noise like Random Noise Attack to get an
 63 optimization variable x' , and the result is shown in figure 4. Since severe noise was generated, we
 64 decided not to use this method.

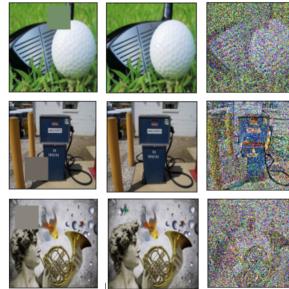


Figure 4: Image inpainting with random samples from pure random noise

65 **Class-conditional seed distribution** Instead, we generated random samples from class-conditional
 66 seed distribution \mathcal{G}_y which was also used in Image Generation, and the result is shown in figure 5.
 67 Additional results can be found in Supplementary Images A.2.

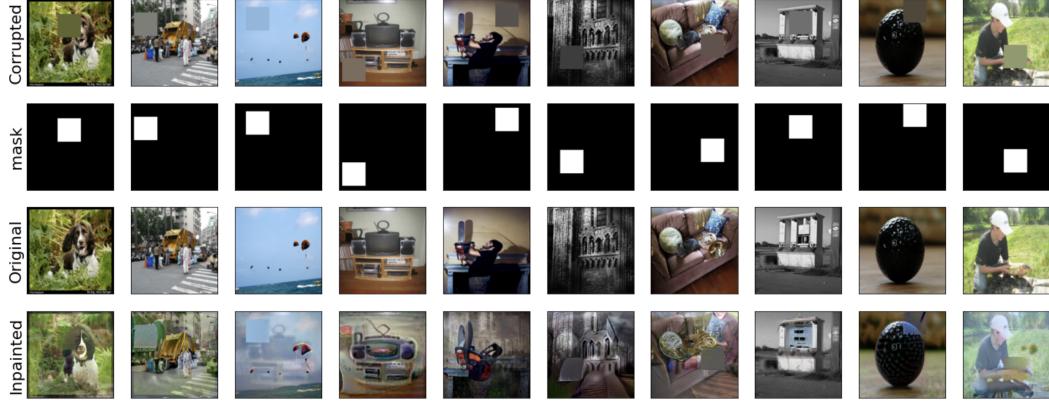


Figure 5: Image inpainting with random samples from class-conditional seed distribution

68 Although inpainted images do not perfectly match the original images, we can say that it's not too
 69 difficult to recognize the images, and most of them look perceptually plausible.

70 4 Super-Resolution

71 The goal of Super Resolution is to take a low-resolution image and increase its resolution. This
 72 is achieved by utilizing a standard L2 FGSM model where we minimize the loss function on an
 73 adversarially trained ImageNet model to increase the resolution of the images. In terms of execution
 74 and evaluation, we took some samples of each class from the Imagenette dataset and down-sampled
 75 the image to reduce its resolution, then we can up-sample this image to restore it to its original 224
 76 x 244 pixel dimensions but still have its low resolution. These images are then passed into the L2
 77 FGSM attack and the resulting output can be seen in the figure below:

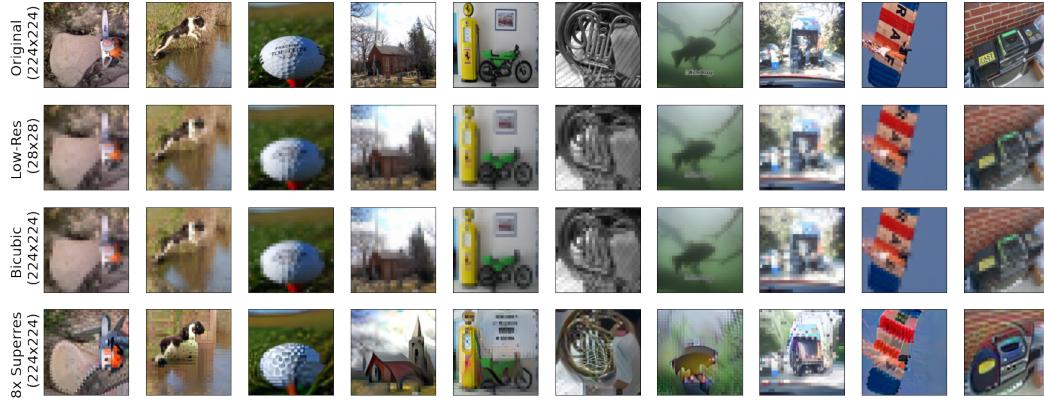


Figure 6: Super-Resolution of images from each class in Imagenette

78 Looking at the output from the L2 FGSM we can see that this methodology did indeed increase the
 79 overall resolution of the object that is classified but there are a few peculiarities to note. The first
 80 interesting point is that in the first column we see that the chainsaw's resolution has been increased
 81 and becomes clear but this method of attack also increases the resolution of the tree in the background
 82 but it is recognized as a saw by the model and it now gains a serrated edge akin to that of a saw-blade.
 83 Moving on to the 6th column which contains a french horn, we notice that our method causes the

84 original black-and-white image to become a colored image where the french horn is now seen as
 85 gold. Lastly in the 7th column with a tench fish, we see that the resolution of the fish itself is greatly
 86 increased but the background becomes distorted. This suggests that the original model utilized the
 87 fish to identify its class but the rest of the image is ignored and thus becomes this blurry and distorted
 88 background. The only scenario in this subset that did not have a clear super-resolution output was
 89 that of the parachute, this can be understood as this parachute seems to be on fire which is potentially
 90 an outlier compared to a standard parachute image.

91 **5 Image-to-Image Translation**

92 For image-to-image translation, our primary goal is to utilize the standard L2 FGSM where the loss
 93 function is minimized to convert Imagenette images from one class to another class. To put it simply
 94 our method of achieving this goal is to take the standard L2 FGSM model with the input images
 95 being the original class and the labels being the class we want to translate the original images into.
 96 Doing this results in the following situations where we translate cassettes to french horns, churches
 97 into gas stations, and garbage trucks into golf balls:

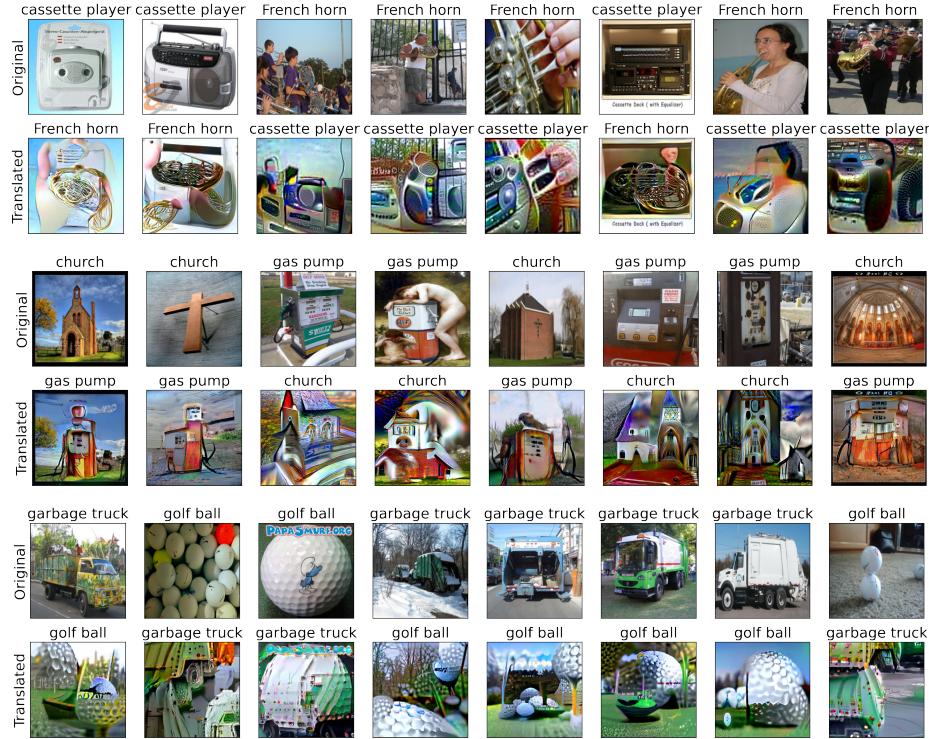


Figure 7: Super Resolution of images from each class in Imagenette

98 From the examples above we can see that this method takes what the model identifies as the original
 99 class and converts those pixels into what it thinks the target class looks like in this scenario. From
 100 the examples, in figure 7 we can see that there is an easily identifiable conversion in each of the
 101 examples. Although the target class is recognizable in the output images, the background of the
 102 image in most cases is either distorted or is replaced with more instances of the target class. This is
 103 easily apparent in the garbage truck to golf ball translation as a singular garbage truck is replaced
 104 with many instances of golf balls. Overall image to image translation in the Imagenette dataset seems
 105 to work but sometimes the outcome is bizarre and wild in some cases.

106 **6 Conclusion**

107 In this work, we leveraged the pre-trained single robustness classifier to perform four different image
 108 synthesis tasks based on the paper released by Santurkar et al (2019).

109 First, in Image Generation, we generated realistic images by adding relevant and significantly
110 meaningful features of the desired class to the input images. Second, in Image Inpainting, we
111 recovered the missing pixels using robust models in a manner that is perceptually plausible for the
112 rest of the image. Third, in Super-Resolution, we took the low-resolution images and increase their
113 resolutions. And the last, in Image-to-Image Translation, we translated images from a source to a
114 target domain in a semantic manner.

115 We were able to obtain results almost similar to those of the paper, and we confirmed the simple
116 feed-forward classifier, when robustly trained, can be a powerful tool for image synthesis.

117 **References**

- 118 [1] Santurkar, S., Tsipras, D., Tran, B., Ilyas, A., Engstrom, L., Madry, A. (2019, August 8). Image synthesis
119 with a single (robust) classifier. arXiv.org. Retrieved December 15, 2022, from <https://arxiv.org/abs/1906.09453>
- 120 [2] Santurkar, S., Tsipras, D., Tran, B., Ilyas, A., Engstrom, L., Madry, A. (2019, August 8) from
121 https://github.com/MadryLab/robustness_applications
- 122 [3] Imagenette dataset by Jeremy Howard from <https://github.com/fastai/imagenette>

123 **A Supplementary Images**

124 **A.1 Image Generation**

125 We generated images from randomly-sampled seeds for 10 classes: tench, English springer, cassette player,
 126 chain saw, church, French horn, garbage truck, gas pump, golf ball, and parachute. The full set of images is
 127 shown in Figure 8. In comparison, the full set of images generated using a standard non-robust classifier is
 128 shown in Figure 9.

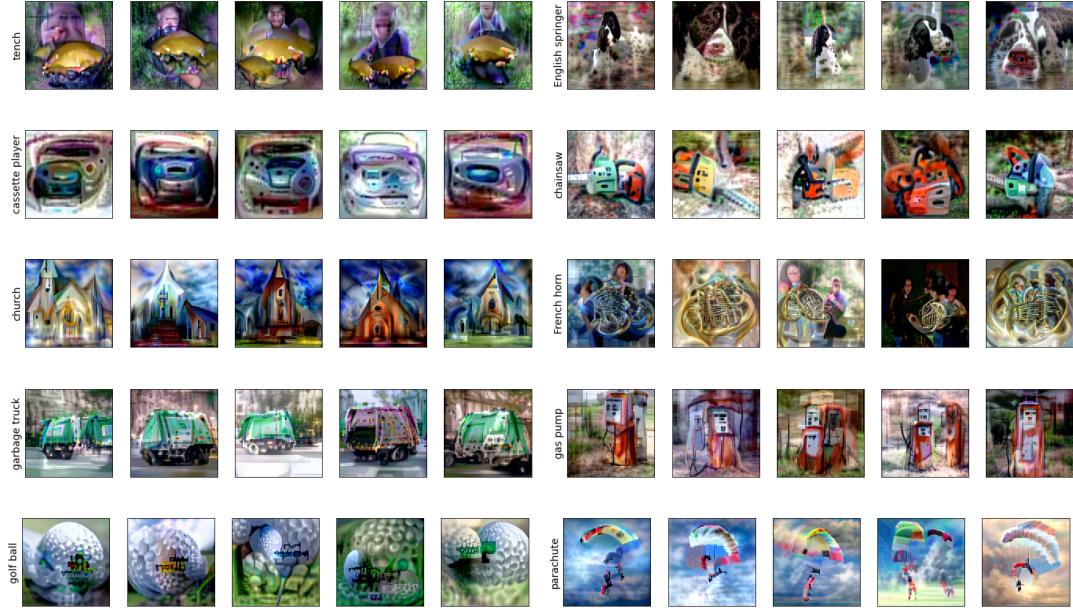


Figure 8: Random samples produced using an adversarially trained robust classifier



Figure 9: Random samples produced using an adversarially trained robust classifier

129 A.2 Image Inpainting

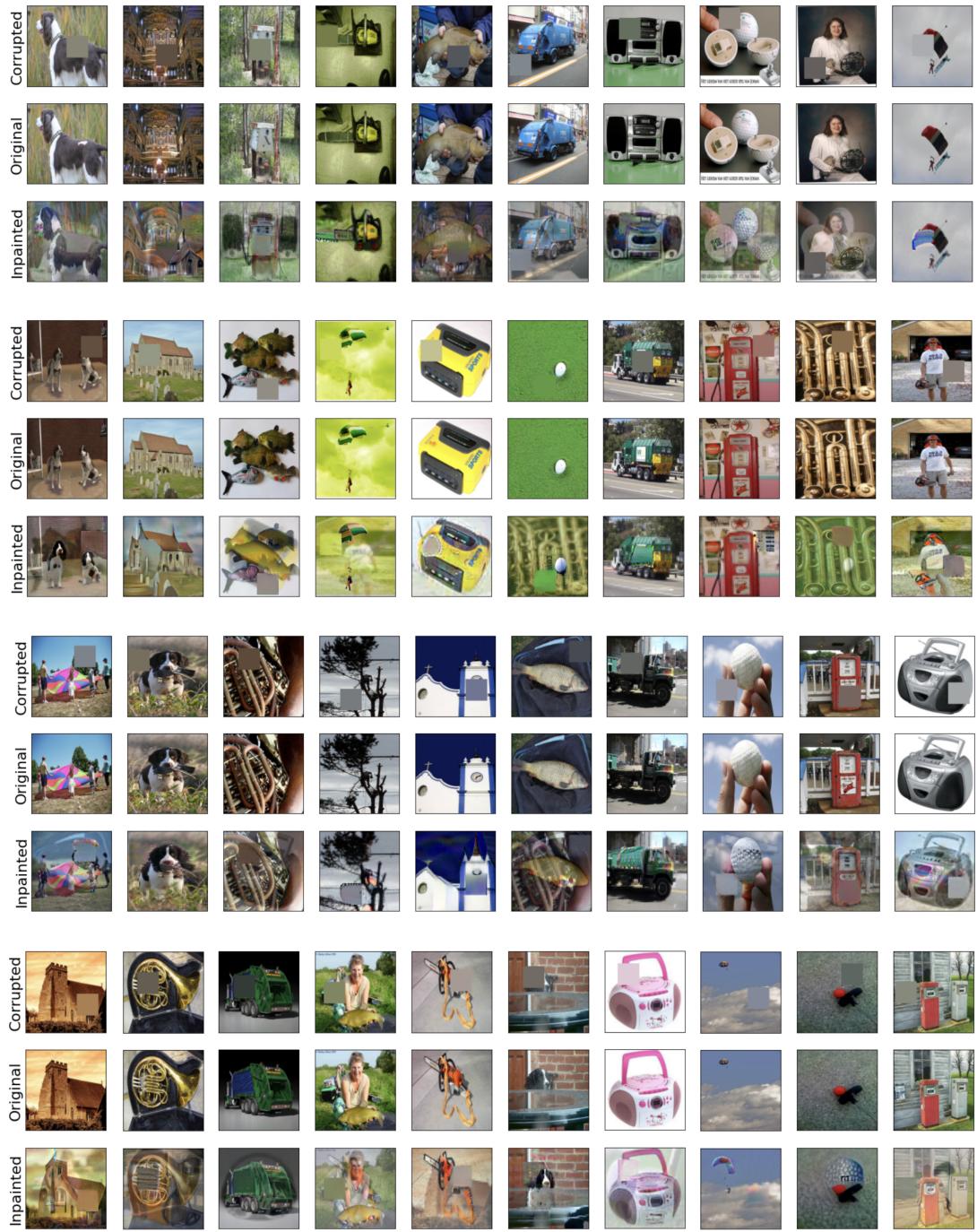


Figure 10: Additional image inpaintings with random samples from class-conditional seed distribution

130 **B Individual Contributions**

- 131 • Shining Yang: Shining is in charge of downloading data, performing the image generation task, making
132 part of the poster, and writing up part of the report.
- 133 • Minjung Lee: Minjung is in charge of setting the project environment, performing the image inpainting,
134 making part of the poster, and writing up part of the report.
- 135 • Abhijith Tammanagari: Abhijith worked Super resolution, image to image translation, making part of
136 the poster, and writing up part of the report.

137 **C Github Repo**

138 Our codes could be found here: ECE661_final_project_f3