

**Implementing the  
Real-Enhanced Super-Resolution Generative Adversarial Network (Real-ESRGAN)**

Georgios Ioannou

Department of Computer Science, The City College Of New York

gioanno000@citymail.cuny.edu

ID: 23927106

<https://github.com/GeorgiosIoannouCoder/realesrgan>

**Abstract.** In the pursuit of enhancing blind super-resolution and upscaling techniques for the restoration of low-resolution images suffered by unknown and intricate degradations, many attempts have been made. However, these attempts have yet to sufficiently address the challenges posed by general real-world degraded images. The paper “Real-ESRGAN: Training Real-World Blind Super-Resolution with Pure Synthetic Data” [1] published on 17, August 2021 introduced an extension of the already existing and successful ESRGAN [2], tailored for practical restoration applications. The training of the Real-ESRGAN involved the utilization of exclusively synthetic data, allowing for a more controlled learning environment. To better simulate the complexities of real-world degradations, a high-order degradation modeling process was incorporated. Furthermore, attention was given to mitigating common artifacts such as ringing and overshoot during the synthesis process. The integration of the U-Net as the discriminator with spectral normalization was a breakthrough as it enhanced the capability of the discriminator used during the training of the Real-ESRGAN and stabilized the training dynamics. Extensive comparisons across various real datasets demonstrated the superiority of the Real-ESRGAN over previous attempts. Therefore, considering the remarkable accomplishments and astonishing outcomes presented in the Real-ESRGAN paper [1], it is highly advisable to allocate time and resources to study this paper more in-depth. Undertaking an in-depth study of the Real-ESRGAN paper [1] is not only merited but also stands as a valuable endeavor, given the substantial impact it has made in the field. Furthermore, the investment of time and effort into the detailed study of Real-ESRGAN is accompanied by the worthwhile pursuit of implementing the code to replicate and reproduce the documented results. The practical engagement in coding to recreate the outcomes outlined in the Real-ESRGAN paper [1] serves as a hands-on approach to understanding the intricacies of the proposed techniques. This active involvement not only solidifies comprehension but also provides an avenue for researchers and practitioners to gain firsthand experience in applying the concepts elucidated in the Real-ESRGAN paper [1]. Real-ESRGAN has a wide range of applications, particularly in the field of image and video processing. Applications include image restoration, image super-resolution, artifact removal, degradation simulation, resize and upscale operations, and JPEG compression.

## 1 Introduction

We will start by defining what super-resolution (SR) is. Single-image super-resolution (SR) constitutes an active research domain with the primary objective of reconstructing a high-resolution (HR) image from its low-resolution (LR) counterpart. The initial strides in this field were marked by the introduction of Super-Resolution Convolutional Neural Network (SRCNN), showcasing the strength of deep convolutional neural network (CNN) approaches and thus leading to significant advancements in SR. Previous blind super-resolution attempts to restore LR images suffered by unknown and intricate degradations. Degradations refer to the processes that transform a high-resolution (HR) image into a low-resolution (LR) version. These degradations are typically used to create training data for SR models. However, the degradations used in training may not perfectly match the degradations experienced in real-world scenarios, which can lead to poor performance of the SR model. Some common degradations used in training SR models include bicubic downsampling and convolution with blur kernels.

Existing strategies can be broadly classified according to the underlying degradation process. The first strategy is called explicit modeling and the second strategy is called implicit modeling. Explicit modeling found in classical degradation models utilizes blur, downsampling, noise, and JPEG compression. Implicit modeling leverages data distribution learning with Generative Adversarial Network (GAN) to discern the degradation model. While effective within training datasets, these methods encounter limitations in generalizing to out-of-distribution images.

In this project, we want to improve the effectiveness of ESRGAN [2] in enhancing regular low-resolution images from the real world. To do this, we will create training examples that mimic common ways images get worse in real life. Actual image problems often arise from a mix of different issues, like how cameras work, changes made to images, and sending pictures online. For instance, when we take photos with our phones, the images can have various issues like blurriness, sensor noise, artificial sharpening effects, and compression due to saving as JPEG files. After that, if we edit the photos and share them on a social media app, more compression and unexpected problems may occur.

These complexities require making the basic “first-level” breakdown model more advanced by considering “higher-level” and more precisely “second-level” breakdowns in real-world situations. This means looking at multiple instances of breakdowns happening and each instance follows the typical breakdown pattern. Moreover, sinc filters are used in the creation process to imitate the usual ringing and overshooting issues. Since Real-ESRGAN is an extension of the ESRGAN the degradation space is much larger which makes the training process difficult and the code more complex. The discriminator needs a more powerful capability to distinguish between real and complex training outputs. The gradient feedback from the discriminator needs to be more accurate for enhancing local details. Therefore, the VGG-style discriminator used in ESRGAN needs to be advanced to a U-Net design. However, the U-Net design structure and complex degradations increase the training instability. Thus, to stabilize the training dynamics, spectral normalization (SN) regularization comes to the rescue. Let us now

see how the new high-order degradation process to model practical degradations, the utilization of sinc filters to model common ringing and overshoot artifacts, and the introduction of the U-Net design create a breakthrough in the field of super-resolution.

## 2 Related Work

One of the first deep learning techniques that has been successfully used for super-resolution was the Super-Resolution Convolutional Neural Network (SRCNN), proposed by Dong and others [3,4]. This method was influential in solving problems such as face hallucination and depth map super-resolution. The SRCNN directly learns a mapping between low-resolution (LR) and high-resolution (HR) images, which makes it faster and more accurate. This is different from other learning-based methods, which need to process the bicubic-upscaled LR images. Wang and others [5] built on this work by replacing the mapping layer with a set of sparse coding sub-networks, creating a sparse coding-based network (SCN). This network outperforms SRCNN with a smaller model size, but it is hard to shrink the sparse coding sub-network without losing mapping accuracy. The Real-ESRGAN does not just work on the original LR image but also contains a simpler but more efficient mapping layer. This network only requires a different deconvolution layer to upscale an image to different sizes, making it faster than other methods that require a completely different network for each upscaling factor.

Generative Adversarial Networks (GANs) are often used to create visually appealing results. They guide the solutions to be closer to the natural image. However, already existing methods use a bicubic downsampling kernel, which often fails when applied to real images. There are two main existing categories of methods for blind super-resolution which both have limitations. The first category uses explicit degradation representations, which typically consist of degradation prediction and conditional restoration. These components can be performed separately or jointly. These methods rely on predefined degradation representations, such as the type and level of degradation, and usually consider simple synthetic degradations. However, if the degradation estimation is not accurate, it can result in artifacts. The second category involves generating training pairs that are as close to real data as possible, and then training a unified network to address blind SR. However, the training pairs are only constrained to degradations associated with specific cameras, and thus may not generalize well to other real images. Moreover, learning fine-grained degradations with unpaired data can be challenging, and the results are often unsatisfactory.

### 3 Second-order Degradation Model

By applying the already existing classical degradation model to generate training pairs, the trained model can handle some real-world samples. However, it still struggles to manage certain complex degradations, particularly unidentified noises and complex artifacts. This is due to the significant difference between synthetic low-resolution images and realistic degraded images. As a result, there is an immediate need to expand the classical degradation model to a high-order model and more specifically to a second-order model to better represent practical degradations. The classical degradation model is essentially a first-order modeling, as it includes a fixed number of basic degradations. However, real-world degradation processes are diverse and often involve a sequence of procedures such as camera imaging systems, image editing, and Internet transmission. For example, restoring a low-quality image downloaded from the Internet involves a complex combination of various degradation processes. The original image, taken with a cellphone many years ago, may contain degradations such as camera blur, sensor noise, low resolution, and JPEG compression. The image is then edited with sharpening and resizing operations, however, this introduces overshoot and blur artifacts. Digital transmission also introduces artifacts, making the process more complicated when the image is shared multiple times on the Internet.

Overshoot artifacts in super-resolution are typically combined with ringing artifacts and are usually produced by sharpening algorithms or JPEG compression. They manifest as an increased jump at the edge transition in an image, causing the image to appear sharper than it actually is.

Blur artifacts in super-resolution can be produced by various resize operations such as nearest-neighbor interpolation, area resize, bilinear interpolation, and bicubic interpolation. These artifacts can result in blurry images when the image is downsampled and then upsampled to its original size.

Therefore, to eliminate overshoot artifacts and blur artifacts there is an urgent need for a second-order degradation model as shown in Equation 1 where each D represents the classical degradation model, y is the input, and x is the output. To further assist the process of eliminating overshoot artifacts and blur artifacts sinc filters are used in the creation process to imitate the usual overshooting and blur issues. Figure 1 shows the overall process of the second-order degradation model used in the Real-EASRGAN training.

$$x = D^2(y) = (D_2 \circ D_1)(y)$$

Equation 1. Second-order Degradation Model Equation

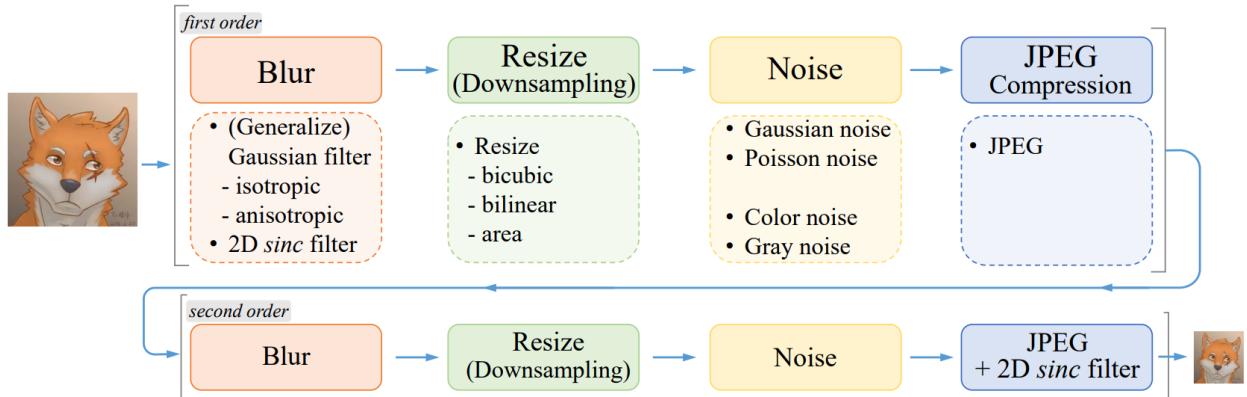


Figure 1. Second-order Degradation Model Process (2-step process)

#### 4. Real-ESRGAN Model Architecture

In the context of super-resolution, the generator and discriminator are key components of a Generative Adversarial Network (GAN) [6].

The generator's role is to produce plausible data, in this case, high-resolution images. The generated instances are used as negative training examples for the discriminator. The generator's goal is to produce output that can fool the discriminator, making it harder for the discriminator to distinguish between real and fake data. As the generator's training progresses, it gets closer to producing output that can fool the discriminator.

On the other hand, the discriminator is a classifier that tries to distinguish real data from fake data created by the generator. It uses real data instances as positive examples and the fake data instances created by the generator as negative examples during training. The discriminator penalizes the generator for producing implausible results, i.e., for producing fake data that it can easily distinguish from real data. The discriminator's goal is to improve its ability to distinguish between real and fake data, which in turn makes it harder for the generator to fool it.

Both the generator and the discriminator are neural networks and require separate training. The generator's output is connected directly to the discriminator's input. Through backpropagation, the discriminator's classification provides a signal that the generator uses to update its weights. In our case, Real-ESRNET is the generator, and U-Net inspired by this paper [7] is the discriminator.

For this paper, we did not train the ESRGAN generator because this will require a total of three trainings instead of two as we will see after. Instead, for the Real-ESRGAN training, we used the ESRGAN model from this paper [7] as the generator. This helped in making the Real-EASRGAN more robust since the ESRGAN is already an effective and tested model. The ESRGAN model architecture as illustrated in Figure 2 is a deep network with several residual-in-residual dense blocks (RRDB). It is worth noting the pixel unshuffle step to reduce the spatial size and enlarge the channel size before feeding inputs into the ESRGAN model.

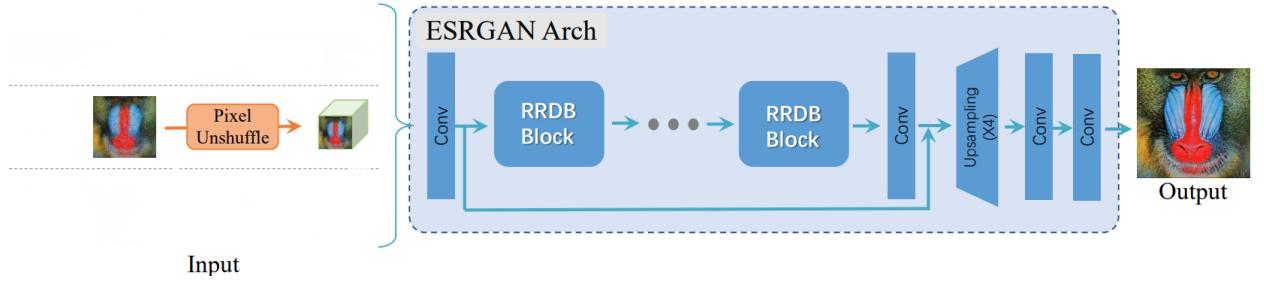


Figure 2. ESRGAN ModelArchitecture

As mentioned before, the training process for the Real-ESRGAN happens in two steps. First, a PSNR-oriented model named Real-ESRNET is trained with the L1 loss using the RRDBNET as the generator and ESRGAN as the base model. In other words, Real-ESRNET is finetuned from ESRGAN. Next, Real-ESRNET is used as an initialization of the generator to train the Real-ESRGAN with a combination of L1 loss, perceptual loss, and GAN loss. If we did not use the ESRGAN model from the paper [7], then we would have to do additional training to get the ESRGAN model.

## 5. Experiments

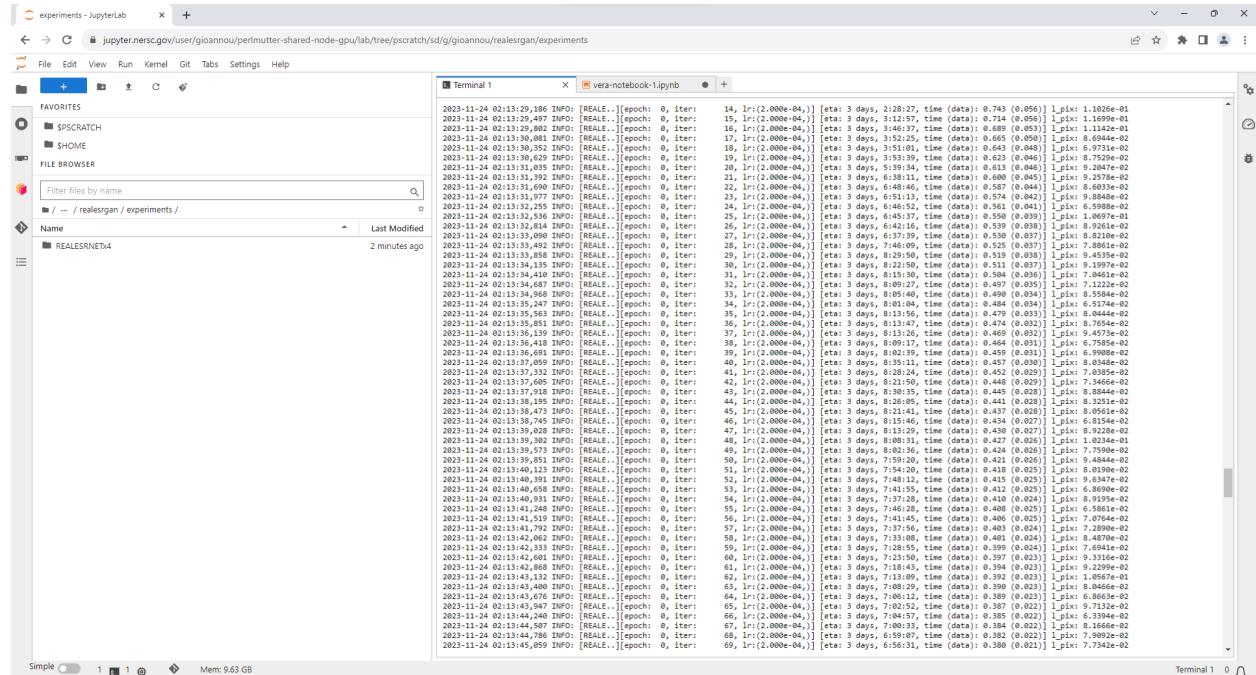
To train the Real-ESRGAN model, we used three datasets. The first dataset is called DIV2K and contains 800 high-resolution (HR) images. The second dataset is called Flickr2K and contains 2,650 high-resolution (HR) images. Both of these datasets have been combined into a single dataset called DF2K and these 3,450 images consist of people, animals, buildings, and scenes. The third dataset is called OutdoorSceneTraining (OST) and is divided into 2,187 high-resolution (HR) images of animals and 2,285 high-resolution (HR) images of buildings. Therefore, the Real-ESRGAN model was trained on a total of 5,735 high-resolution (HR) images.

To train the Real-ESRNET and the Real-ESRGAN four NVIDIA A100 80GB PCIe GPUs were utilized through NERSC’s exclusive GPU nodes. The training for the Real-ESRNET took three days and the training for the Real-ESRGAN took two days. Figure 3 and Figure 4 demonstrate the training for the Real-ESRNET and the Real-ESRGAN respectively. The figures clearly show the loss per pixel for each and every iteration of the training.

For both training sessions, the image size for the ground truth images was set to 256 during the training. Batch size was set to 48 and Adam was used as the optimizer. The L1 loss was used for the Real-ESRNET training. A combination of L1 loss, perceptual loss, and GAN loss was used for the Real-ESRGAN training. That is the reason why Figure 4, shows more loss per pixel values for each iteration in the training. LeakyReLU was used as the activation function. Real-ESRNET was trained with 1,000,000 iterations and a learning rate of  $2 \times 10^{-4}$ . Real-ESRGAN was trained with 400,000 iterations and a learning rate of  $1 \times 10^{-4}$ . All the

training parameters can be found in the .yml files inside the training\_parameters directory of this project's GitHub repository.

Furthermore, Gaussian kernels, generalized Gaussian kernels, and plateau-shaped kernels were utilized, each assigned with a probability of 0.7, 0.15, and 0.15 respectively. The size of the blur kernel is chosen randomly from (7, 9, 11, 13, 15, 17, 19, 21). The blur standard deviation,  $\sigma$ , is sampled from the range [0.2, 3] for the first degradation process and from the range [0.2, 1.5] for the second degradation process. The shape parameter  $\beta$  is sampled from [0.5, 4] and [1, 2] for generalized Gaussian and plateau-shaped kernels, respectively. Additionally, a sinc kernel is introduced with a 0.1 probability, and there is a 0.2 probability of skipping the second blur degradation. For noise, Gaussian and Poisson noises are used with probabilities of {0.5, 0.5}. The noise sigma range and Poisson noise scale are specified as [1, 30] and [0.05, 3], respectively, and [1, 25] and [0.05, 2.5] for the second degradation process. The probability of gray noise is set to 0.4. JPEG compression quality factor ranges from 30 to 95. Finally, the application of the sinc filter is determined with a probability of 0.8.



The screenshot shows a Jupyter Notebook environment with a terminal window displaying training logs for a Real-ESRNet model. The terminal output is as follows:

```

2023-11-24 02:13:29,186 INFO: [REALE].[epoch: 0, iter: 14, lr:(2.000e-04)] etat: 3 days, 2:28:27, time (data): 0.743 (0.056) l_pix: 1.1026e-01
2023-11-24 02:13:29,497 INFO: [REALE].[epoch: 0, iter: 15, lr:(2.000e-04)] etat: 3 days, 3:12:57, time (data): 0.714 (0.056) l_pix: 1.1699e-01
2023-11-24 02:13:29,613 INFO: [REALE].[epoch: 0, iter: 16, lr:(2.000e-04)] etat: 3 days, 3:18:37, time (data): 0.693 (0.053) l_pix: 1.5424e-01
2023-11-24 02:13:30,043 INFO: [REALE].[epoch: 0, iter: 17, lr:(2.000e-04)] etat: 3 days, 3:25:27, time (data): 0.665 (0.052) l_pix: 1.6000e-01
2023-11-24 02:13:30,352 INFO: [REALE].[epoch: 0, iter: 18, lr:(2.000e-04)] etat: 3 days, 3:51:01, time (data): 0.645 (0.048) l_pix: 0.9731e-02
2023-11-24 02:13:30,629 INFO: [REALE].[epoch: 0, iter: 19, lr:(2.000e-04)] etat: 3 days, 3:53:39, time (data): 0.623 (0.046) l_pix: 0.7529e-02
2023-11-24 02:13:31,007 INFO: [REALE].[epoch: 0, iter: 20, lr:(2.000e-04)] etat: 3 days, 3:56:11, time (data): 0.602 (0.045) l_pix: 0.6000e-02
2023-11-24 02:13:31,392 INFO: [REALE].[epoch: 0, iter: 21, lr:(2.000e-04)] etat: 3 days, 6:38:11, time (data): 0.600 (0.045) l_pix: 0.2578e-02
2023-11-24 02:13:31,690 INFO: [REALE].[epoch: 0, iter: 22, lr:(2.000e-04)] etat: 3 days, 6:48:46, time (data): 0.587 (0.044) l_pix: 0.6033e-02
2023-11-24 02:13:32,078 INFO: [REALE].[epoch: 0, iter: 23, lr:(2.000e-04)] etat: 3 days, 6:55:11, time (data): 0.570 (0.040) l_pix: 0.4000e-02
2023-11-24 02:13:32,255 INFO: [REALE].[epoch: 0, iter: 24, lr:(2.000e-04)] etat: 3 days, 6:46:52, time (data): 0.561 (0.041) l_pix: 0.5988e-02
2023-11-24 02:13:32,492 INFO: [REALE].[epoch: 0, iter: 25, lr:(2.000e-04)] etat: 3 days, 6:45:37, time (data): 0.559 (0.039) l_pix: 1.0097e-01
2023-11-24 02:13:33,043 INFO: [REALE].[epoch: 0, iter: 26, lr:(2.000e-04)] etat: 3 days, 6:44:22, time (data): 0.557 (0.039) l_pix: 0.6000e-02
2023-11-24 02:13:33,090 INFO: [REALE].[epoch: 0, iter: 27, lr:(2.000e-04)] etat: 3 days, 6:37:39, time (data): 0.550 (0.037) l_pix: 0.8230e-02
2023-11-24 02:13:33,492 INFO: [REALE].[epoch: 0, iter: 28, lr:(2.000e-04)] etat: 3 days, 7:46:09, time (data): 0.525 (0.037) l_pix: 7.8861e-02
2023-11-24 02:13:33,536 INFO: [REALE].[epoch: 0, iter: 29, lr:(2.000e-04)] etat: 3 days, 8:29:50, time (data): 0.515 (0.038) l_pix: 0.4035e-02
2023-11-24 02:13:33,538 INFO: [REALE].[epoch: 0, iter: 30, lr:(2.000e-04)] etat: 3 days, 8:30:27, time (data): 0.511 (0.038) l_pix: 0.4000e-02
2023-11-24 02:13:34,410 INFO: [REALE].[epoch: 0, iter: 31, lr:(2.000e-04)] etat: 3 days, 8:15:30, time (data): 0.508 (0.036) l_pix: 7.0461e-02
2023-11-24 02:13:34,487 INFO: [REALE].[epoch: 0, iter: 32, lr:(2.000e-04)] etat: 3 days, 8:09:27, time (data): 0.497 (0.035) l_pix: 7.1222e-02
2023-11-24 02:13:34,520 INFO: [REALE].[epoch: 0, iter: 33, lr:(2.000e-04)] etat: 3 days, 8:08:54, time (data): 0.495 (0.035) l_pix: 0.4000e-02
2023-11-24 02:13:35,247 INFO: [REALE].[epoch: 0, iter: 34, lr:(2.000e-04)] etat: 3 days, 8:01:04, time (data): 0.484 (0.034) l_pix: 0.5174e-02
2023-11-24 02:13:35,563 INFO: [REALE].[epoch: 0, iter: 35, lr:(2.000e-04)] etat: 3 days, 8:13:56, time (data): 0.479 (0.033) l_pix: 0.4044e-02
2023-11-24 02:13:35,613 INFO: [REALE].[epoch: 0, iter: 36, lr:(2.000e-04)] etat: 3 days, 8:14:21, time (data): 0.477 (0.033) l_pix: 0.4000e-02
2023-11-24 02:13:36,139 INFO: [REALE].[epoch: 0, iter: 37, lr:(2.000e-04)] etat: 3 days, 8:13:26, time (data): 0.469 (0.032) l_pix: 0.4573e-02
2023-11-24 02:13:36,418 INFO: [REALE].[epoch: 0, iter: 38, lr:(2.000e-04)] etat: 3 days, 8:09:17, time (data): 0.466 (0.031) l_pix: 0.7585e-02
2023-11-24 02:13:36,448 INFO: [REALE].[epoch: 0, iter: 39, lr:(2.000e-04)] etat: 3 days, 8:09:39, time (data): 0.464 (0.031) l_pix: 0.4000e-02
2023-11-24 02:13:37,059 INFO: [REALE].[epoch: 0, iter: 40, lr:(2.000e-04)] etat: 3 days, 8:09:33, time (data): 0.457 (0.030) l_pix: 0.8048e-02
2023-11-24 02:13:37,332 INFO: [REALE].[epoch: 0, iter: 41, lr:(2.000e-04)] etat: 3 days, 8:28:24, time (data): 0.452 (0.029) l_pix: 7.0385e-02
2023-11-24 02:13:37,777 INFO: [REALE].[epoch: 0, iter: 42, lr:(2.000e-04)] etat: 3 days, 8:21:50, time (data): 0.446 (0.029) l_pix: 7.0366e-02
2023-11-24 02:13:38,119 INFO: [REALE].[epoch: 0, iter: 43, lr:(2.000e-04)] etat: 3 days, 8:20:27, time (data): 0.443 (0.028) l_pix: 0.8000e-02
2023-11-24 02:13:38,195 INFO: [REALE].[epoch: 0, iter: 44, lr:(2.000e-04)] etat: 3 days, 8:26:05, time (data): 0.441 (0.028) l_pix: 0.3251e-02
2023-11-24 02:13:38,473 INFO: [REALE].[epoch: 0, iter: 45, lr:(2.000e-04)] etat: 3 days, 8:21:41, time (data): 0.437 (0.028) l_pix: 0.0561e-02
2023-11-24 02:13:38,503 INFO: [REALE].[epoch: 0, iter: 46, lr:(2.000e-04)] etat: 3 days, 8:21:46, time (data): 0.435 (0.028) l_pix: 0.4000e-02
2023-11-24 02:13:39,028 INFO: [REALE].[epoch: 0, iter: 47, lr:(2.000e-04)] etat: 3 days, 8:13:29, time (data): 0.430 (0.027) l_pix: 0.9228e-02
2023-11-24 02:13:39,472 INFO: [REALE].[epoch: 0, iter: 48, lr:(2.000e-04)] etat: 3 days, 8:00:31, time (data): 0.427 (0.026) l_pix: 1.0234e-01
2023-11-24 02:13:39,502 INFO: [REALE].[epoch: 0, iter: 49, lr:(2.000e-04)] etat: 3 days, 7:59:50, time (data): 0.425 (0.026) l_pix: 0.4000e-02
2023-11-24 02:13:39,851 INFO: [REALE].[epoch: 0, iter: 50, lr:(2.000e-04)] etat: 3 days, 7:59:20, time (data): 0.421 (0.026) l_pix: 0.9444e-02
2023-11-24 02:13:40,123 INFO: [REALE].[epoch: 0, iter: 51, lr:(2.000e-04)] etat: 3 days, 7:54:20, time (data): 0.411 (0.025) l_pix: 0.0190e-02
2023-11-24 02:13:40,143 INFO: [REALE].[epoch: 0, iter: 52, lr:(2.000e-04)] etat: 3 days, 7:54:40, time (data): 0.410 (0.025) l_pix: 0.4000e-02
2023-11-24 02:13:40,480 INFO: [REALE].[epoch: 0, iter: 53, lr:(2.000e-04)] etat: 3 days, 7:41:55, time (data): 0.412 (0.025) l_pix: 0.8690e-02
2023-11-24 02:13:40,931 INFO: [REALE].[epoch: 0, iter: 54, lr:(2.000e-04)] etat: 3 days, 7:37:28, time (data): 0.410 (0.024) l_pix: 0.9195e-02
2023-11-24 02:13:41,248 INFO: [REALE].[epoch: 0, iter: 55, lr:(2.000e-04)] etat: 3 days, 7:46:28, time (data): 0.408 (0.025) l_pix: 0.5001e-02
2023-11-24 02:13:41,441 INFO: [REALE].[epoch: 0, iter: 56, lr:(2.000e-04)] etat: 3 days, 7:41:45, time (data): 0.407 (0.025) l_pix: 0.4000e-02
2023-11-24 02:13:41,792 INFO: [REALE].[epoch: 0, iter: 57, lr:(2.000e-04)] etat: 3 days, 7:37:56, time (data): 0.403 (0.024) l_pix: 7.2890e-02
2023-11-24 02:13:42,062 INFO: [REALE].[epoch: 0, iter: 58, lr:(2.000e-04)] etat: 3 days, 7:33:08, time (data): 0.401 (0.024) l_pix: 0.4070e-02
2023-11-24 02:13:42,332 INFO: [REALE].[epoch: 0, iter: 59, lr:(2.000e-04)] etat: 3 days, 7:30:28, time (data): 0.399 (0.023) l_pix: 0.4000e-02
2023-11-24 02:13:42,601 INFO: [REALE].[epoch: 0, iter: 60, lr:(2.000e-04)] etat: 3 days, 7:23:50, time (data): 0.397 (0.023) l_pix: 0.9311e-02
2023-11-24 02:13:42,868 INFO: [REALE].[epoch: 0, iter: 61, lr:(2.000e-04)] etat: 3 days, 7:18:43, time (data): 0.394 (0.023) l_pix: 0.2299e-02
2023-11-24 02:13:43,138 INFO: [REALE].[epoch: 0, iter: 62, lr:(2.000e-04)] etat: 3 days, 7:14:13, time (data): 0.392 (0.023) l_pix: 0.4000e-02
2023-11-24 02:13:43,400 INFO: [REALE].[epoch: 0, iter: 63, lr:(2.000e-04)] etat: 3 days, 7:08:29, time (data): 0.390 (0.023) l_pix: 0.8466e-02
2023-11-24 02:13:43,676 INFO: [REALE].[epoch: 0, iter: 64, lr:(2.000e-04)] etat: 3 days, 7:06:12, time (data): 0.388 (0.023) l_pix: 0.8663e-02
2023-11-24 02:13:44,044 INFO: [REALE].[epoch: 0, iter: 65, lr:(2.000e-04)] etat: 3 days, 7:04:57, time (data): 0.386 (0.022) l_pix: 0.8339e-02
2023-11-24 02:13:44,240 INFO: [REALE].[epoch: 0, iter: 66, lr:(2.000e-04)] etat: 3 days, 7:03:30, time (data): 0.385 (0.022) l_pix: 0.8339e-02
2023-11-24 02:13:44,507 INFO: [REALE].[epoch: 0, iter: 67, lr:(2.000e-04)] etat: 3 days, 7:00:33, time (data): 0.384 (0.022) l_pix: 0.1666e-02
2023-11-24 02:13:44,780 INFO: [REALE].[epoch: 0, iter: 68, lr:(2.000e-04)] etat: 3 days, 6:59:07, time (data): 0.382 (0.022) l_pix: 7.9002e-02
2023-11-24 02:13:45,059 INFO: [REALE].[epoch: 0, iter: 69, lr:(2.000e-04)] etat: 3 days, 6:56:31, time (data): 0.380 (0.021) l_pix: 7.7342e-02

```

Figure 3. Real-ESRNET Training

```

t_d_realm: 3.8790e-01 l_d_fake: 6.4919e-01 out_d_fake: -3.6542e-02
2023-11-24 02:17:30,551 INFO: [REAL-ESRGAN]:[Epoch: 0, iter: 85, lr:(1.000e-04)] [eta: 2 days, 1:35:23, time (data): 0.597 (0.030)] l_g_pix: 5.9443e-02 l_g_percep: 1.3766e+01 l_g_gen: 7.9321e-02 l_d_real: 7.3404e-01 ou
t_d_realm: -5.0649e-02 l_d_fake: 6.0821e-01 out_d_fake: -1.8489e-01
2023-11-24 02:17:30,579 INFO: [REAL-ESRGAN]:[Epoch: 0, iter: 86, lr:(1.000e-04)] [eta: 2 days, 1:33:56, time (data): 0.596 (0.029)] l_g_pix: 5.4089e-02 l_g_percep: 1.0471e+01 l_g_gen: 7.9102e-02 l_d_real: 7.2898e-01 ou
t_d_realm: -5.0649e-02 l_d_fake: 6.0821e-01 out_d_fake: -1.8377e-01
2023-11-24 02:17:31,402 INFO: [REAL-ESRGAN]:[Epoch: 0, iter: 87, lr:(1.000e-04)] [eta: 2 days, 1:32:14, time (data): 0.594 (0.029)] l_g_pix: 6.9411e-02 l_g_percep: 1.4107e+01 l_g_gen: 7.8721e-02 l_d_real: 6.7390e-01 ou
t_d_realm: 1.2104e-01 l_d_fake: 6.2042e-01 out_d_fake: -1.6657e-01
2023-11-24 02:17:31,420 INFO: [REAL-ESRGAN]:[Epoch: 0, iter: 88, lr:(1.000e-04)] [eta: 2 days, 1:31:09, time (data): 0.593 (0.029)] l_g_pix: 6.2042e-01 l_g_percep: 9.2014e+00 l_g_gen: 7.8591e-02 l_d_real: 5.0433e-01 ou
t_d_realm: 6.6401e-01 l_d_fake: 7.3342e-01 out_d_fake: -6.6380e-02
2023-11-24 02:17:31,428 INFO: [REAL-ESRGAN]:[Epoch: 0, iter: 89, lr:(1.000e-04)] [eta: 2 days, 1:33:40, time (data): 0.596 (0.029)] l_g_pix: 6.2460e-02 l_g_percep: 1.2214e+01 l_g_gen: 7.5951e-02 l_d_real: 6.5490e-01 ou
t_d_realm: 4.8030e-01 l_d_fake: 7.0302e-01 out_d_fake: -1.0100e-01
2023-11-24 02:17:31,446 INFO: [REAL-ESRGAN]:[Epoch: 0, iter: 90, lr:(1.000e-04)] [eta: 2 days, 1:32:46, time (data): 0.589 (0.028)] l_g_pix: 9.7567e-02 l_g_percep: 1.6613e+01 l_g_gen: 7.2089e-02 l_d_real: 5.6028e-01 ou
t_d_realm: 4.0417e-01 l_d_fake: 6.7418e-02 out_d_fake: -2.0631e-01
2023-11-24 02:17:31,464 INFO: [REAL-ESRGAN]:[Epoch: 0, iter: 91, lr:(1.000e-04)] [eta: 2 days, 1:29:31, time (data): 0.585 (0.028)] l_g_pix: 7.0706e-02 l_g_percep: 1.4442e+01 l_g_gen: 8.8845e-02 l_d_real: 7.6447e-01 ou
t_d_realm: 3.9146e-01 l_d_fake: 5.8668e-01 out_d_fake: -3.9215e-01
2023-11-24 02:17:31,482 INFO: [REAL-ESRGAN]:[Epoch: 0, iter: 92, lr:(1.000e-04)] [eta: 2 days, 1:28:05, time (data): 0.585 (0.028)] l_g_pix: 7.1713e-02 l_g_percep: 1.5095e+01 l_g_gen: 8.4197e-02 l_d_real: 6.4642e-01 ou
t_d_realm: 3.9146e-01 l_d_fake: 5.8668e-01 out_d_fake: -2.5449e-01
2023-11-24 02:17:31,499 INFO: [REAL-ESRGAN]:[Epoch: 0, iter: 93, lr:(1.000e-04)] [eta: 2 days, 1:26:54, time (data): 0.582 (0.027)] l_g_pix: 8.7315e-02 l_g_percep: 1.8343e+01 l_g_gen: 5.9311e-02 l_d_real: 3.8301e-01 ou
t_d_realm: 6.3028e-02 l_d_fake: 5.2358e-01 out_d_fake: -3.8604e-01
2023-11-24 02:17:31,517 INFO: [REAL-ESRGAN]:[Epoch: 0, iter: 94, lr:(1.000e-04)] [eta: 2 days, 1:26:29, time (data): 0.580 (0.027)] l_g_pix: 8.3690e-02 l_g_percep: 1.6127e+01 l_g_gen: 9.1092e-02 l_d_real: 7.3845e-01 ou
t_d_realm: 4.7394e-02 l_d_fake: 4.7394e-01 out_d_fake: -4.7394e-01
2023-11-24 02:17:31,535 INFO: [REAL-ESRGAN]:[Epoch: 0, iter: 95, lr:(1.000e-04)] [eta: 2 days, 1:26:15, time (data): 0.579 (0.027)] l_g_pix: 6.5262e-02 l_g_percep: 1.3709e+01 l_g_gen: 9.6262e-02 l_d_real: 9.1725e-01 ou
t_d_realm: -1.9005e-01 l_d_fake: 4.0712e-01 out_d_fake: -3.1875e-01
2023-11-24 02:17:31,553 INFO: [REAL-ESRGAN]:[Epoch: 0, iter: 96, lr:(1.000e-04)] [eta: 2 days, 1:26:35, time (data): 0.577 (0.027)] l_g_pix: 8.4283e-02 l_g_percep: 1.5392e+01 l_g_gen: 9.3325e-02 l_d_real: 8.6778e-01 ou
t_d_realm: -2.0000e-01 l_d_fake: 5.3585e-01 out_d_fake: -4.2668e-01
2023-11-24 02:17:31,571 INFO: [REAL-ESRGAN]:[Epoch: 0, iter: 97, lr:(1.000e-04)] [eta: 2 days, 1:25:04, time (data): 0.576 (0.027)] l_g_pix: 7.0576e-02 l_g_percep: 1.4442e+01 l_g_gen: 8.8845e-02 l_d_real: 7.6447e-01 ou
t_d_realm: -3.5615e-01 l_d_fake: 5.1573e-01 out_d_fake: -3.9892e-01
2023-11-24 02:17:31,589 INFO: [REAL-ESRGAN]:[Epoch: 0, iter: 98, lr:(1.000e-04)] [eta: 2 days, 1:24:19, time (data): 0.574 (0.027)] l_g_pix: 8.7315e-02 l_g_percep: 1.8343e+01 l_g_gen: 5.9311e-02 l_d_real: 3.8301e-01 ou
t_d_realm: -4.7238e-01 l_d_fake: 5.1573e-01 out_d_fake: -4.7238e-01
2023-11-24 02:17:31,607 INFO: [REAL-ESRGAN]:[Epoch: 0, iter: 99, lr:(1.000e-04)] [eta: 2 days, 1:24:17, time (data): 0.573 (0.027)] l_g_pix: 7.7080e-02 l_g_percep: 1.6126e+01 l_g_gen: 8.7243e-02 l_d_real: 8.4563e-01 ou
t_d_realm: -2.7699e-01 l_d_fake: 5.4516e-01 out_d_fake: -3.2691e-01
2023-11-24 02:17:31,625 INFO: [REAL-ESRGAN]:[Epoch: 0, iter: 100, lr:(1.000e-04)] [eta: 2 days, 1:23:20, time (data): 0.572 (0.026)] l_g_pix: 6.5262e-02 l_g_percep: 1.3709e+01 l_g_gen: 8.6693e-02 l_d_real: 8.4673e-01 ou
t_d_realm: -2.8314e-01 l_d_fake: 5.4795e-01 out_d_fake: -3.1875e-01
2023-11-24 02:17:31,643 INFO: [REAL-ESRGAN]:[Epoch: 0, iter: 101, lr:(1.000e-04)] [eta: 2 days, 1:23:39, time (data): 0.571 (0.026)] l_g_pix: 7.8397e-02 l_g_percep: 1.5666e+01 l_g_gen: 8.4138e-02 l_d_real: 8.1804e-01 ou
t_d_realm: -2.7239e-01 l_d_fake: 5.4795e-01 out_d_fake: -2.7239e-01
2023-11-24 02:17:31,661 INFO: [REAL-ESRGAN]:[Epoch: 0, iter: 102, lr:(1.000e-04)] [eta: 2 days, 1:23:03, time (data): 0.570 (0.026)] l_g_pix: 7.1698e-02 l_g_percep: 1.4458e+01 l_g_gen: 7.9404e-02 l_d_real: 7.8137e-01 ou
t_d_realm: -1.6755e-01 l_d_fake: 6.0397e-01 out_d_fake: -1.9030e-01
2023-11-24 02:17:31,679 INFO: [REAL-ESRGAN]:[Epoch: 0, iter: 103, lr:(1.000e-04)] [eta: 2 days, 1:21:46, time (data): 0.569 (0.026)] l_g_pix: 8.4857e-02 l_g_percep: 1.5756e+01 l_g_gen: 7.4916e-02 l_d_real: 7.3730e-01 ou
t_d_realm: -8.5055e-02 l_d_fake: 5.4157e-01 out_d_fake: -1.0732e-01
2023-11-24 02:17:31,697 INFO: [REAL-ESRGAN]:[Epoch: 0, iter: 104, lr:(1.000e-04)] [eta: 2 days, 1:25:20, time (data): 0.568 (0.026)] l_g_pix: 8.4857e-02 l_g_percep: 1.5756e+01 l_g_gen: 7.4916e-02 l_d_real: 7.3730e-01 ou
t_d_realm: -5.8243e-02 l_d_fake: 6.4597e-01 out_d_fake: -9.9331e-02
2023-11-24 02:17:31,715 INFO: [REAL-ESRGAN]:[Epoch: 0, iter: 105, lr:(1.000e-04)] [eta: 2 days, 1:25:32, time (data): 0.567 (0.026)] l_g_pix: 8.2533e-02 l_g_percep: 1.5726e+01 l_g_gen: 7.4552e-02 l_d_real: 7.2579e-01 ou
t_d_realm: -2.0000e-01 l_d_fake: 6.4597e-01 out_d_fake: -2.0000e-01
2023-11-24 02:17:31,733 INFO: [REAL-ESRGAN]:[Epoch: 0, iter: 106, lr:(1.000e-04)] [eta: 2 days, 1:23:55, time (data): 0.566 (0.026)] l_g_pix: 7.3340e-02 l_g_percep: 1.4421e+01 l_g_gen: 7.2909e-02 l_d_real: 7.0418e-01 ou
t_d_realm: -8.4661e-03 l_d_fake: 6.6126e-01 out_d_fake: -6.8404e-02
2023-11-24 02:17:31,751 INFO: [REAL-ESRGAN]:[Epoch: 0, iter: 107, lr:(1.000e-04)] [eta: 2 days, 1:22:30, time (data): 0.564 (0.026)] l_g_pix: 8.6004e-02 l_g_percep: 1.5455e+01 l_g_gen: 7.0843e-02 l_d_real: 6.8574e-01 ou
t_d_realm: -1.0000e-01 l_d_fake: 6.6126e-01 out_d_fake: -1.0000e-01
2023-11-24 02:17:31,769 INFO: [REAL-ESRGAN]:[Epoch: 0, iter: 108, lr:(1.000e-04)] [eta: 2 days, 1:22:05, time (data): 0.563 (0.025)] l_g_pix: 7.1373e+01 l_g_percep: 1.1644e+01 l_g_gen: 6.8185e-02 l_d_real: 6.6526e-01 ou
t_d_realm: 6.6651e-03 l_d_fake: 7.1122e-01 out_d_fake: 2.9362e-02
2023-11-24 02:17:31,787 INFO: [REAL-ESRGAN]:[Epoch: 0, iter: 109, lr:(1.000e-04)] [eta: 2 days, 1:20:41, time (data): 0.562 (0.025)] l_g_pix: 7.2257e-02 l_g_percep: 1.3282e+01 l_g_gen: 6.6380e-02 l_d_real: 6.3879e-01 ou
t_d_realm: 1.3038e-01 l_d_fake: 7.3086e-01 out_d_fake: 6.6895e-02
2023-11-24 02:17:31,805 INFO: [REAL-ESRGAN]:[Epoch: 0, iter: 110, lr:(1.000e-04)] [eta: 2 days, 1:19:26, time (data): 0.561 (0.025)] l_g_pix: 5.5522e-02 l_g_percep: 1.2320e+01 l_g_gen: 6.5215e-02 l_d_real: 6.3918e-01 ou
t_d_realm: 1.3038e-01 l_d_fake: 7.3086e-01 out_d_fake: 1.0154e-01
2023-11-24 02:17:31,823 INFO: [REAL-ESRGAN]:[Epoch: 0, iter: 111, lr:(1.000e-04)] [eta: 2 days, 1:28:42, time (data): 0.561 (0.025)] l_g_pix: 6.9696e-02 l_g_percep: 1.5142e+01 l_g_gen: 6.5584e-02 l_d_real: 6.2501e-01 ou
t_d_realm: 1.7959e-01 l_d_fake: 7.5765e-01 out_d_fake: 1.0181e-01

```

Figure 4. Real-ESRGAN Training

## 6. Coding Files

All the coding files (highly documented), steps on how to run them, steps on how to structure this project, and detailed video demos can all be found in this project's GitHub repository README.md file. Let us mention the purpose of each file.

The file `unet_discriminator_sn.py` defines the architecture of the U-Net discriminator used for training the Real-ESRGAN. The file `real_esrnet_model.py` defines The architecture of the generator Real-ESRNET used for training the Real-ESRGAN. The file `image_crop.py` crops the images of the DF2K dataset to sub-images for faster IO. The file `image_path.py` generates all the ground-truth image paths to be used during training. The file `multiscale_df2k.py` scales up and down images. The scales that are used are 0.75, 0.5, and  $\frac{1}{3}$ . The files `inference.py` and `real_esrgan.py` serve as an inference for the Real-ESRGAN model. The file `real_esrgan_model.py` defines the architecture of the Real-ESRGAN model as described in this project and this paper [1]. The file `data_used_during_training.py` augments the datasets so that the Real-ESRGAN model can be trained on only pure synthetic data. The file `training.py` is used to train both the Real-ESRNET and the Real-ESRGAN models.

## 7. Differences From The Original Paper

This project differs slightly from the original paper [1] in two ways. The first difference is that this project only trained the Real-ESRGAN for a scale x4, but the original paper also trained the Real-ESRGAN model for the scales x1 and x2. This project chose to train the Real-ESRGAN model for the scale x4 because based on a previous paper [2], the best tradeoff between performance and parameters is the scale factor x4. The second difference is that the original paper [1], compared their Real-ESRGAN model with ESRGAN, DAN, CDC, RealSR, and BSRGAN using the testing datasets RealSR, DRealSR, OST300, DPED, and ADE20K. This project could not do that due to the difficulty of finding these datasets and the models from previous works.

## 8. Compare Results To Original Paper

The results from this project can be found in this project's GitHub repository README.md file as images and video demos. Since this project followed the same methodologies as the original paper [1], the results were expected to be the same and as accurate. Figure 5 and Figure 6, illustrate the results of this project's model on images that are mentioned in the original paper [1]. Figure 7 shows this project's model performance on a building image.

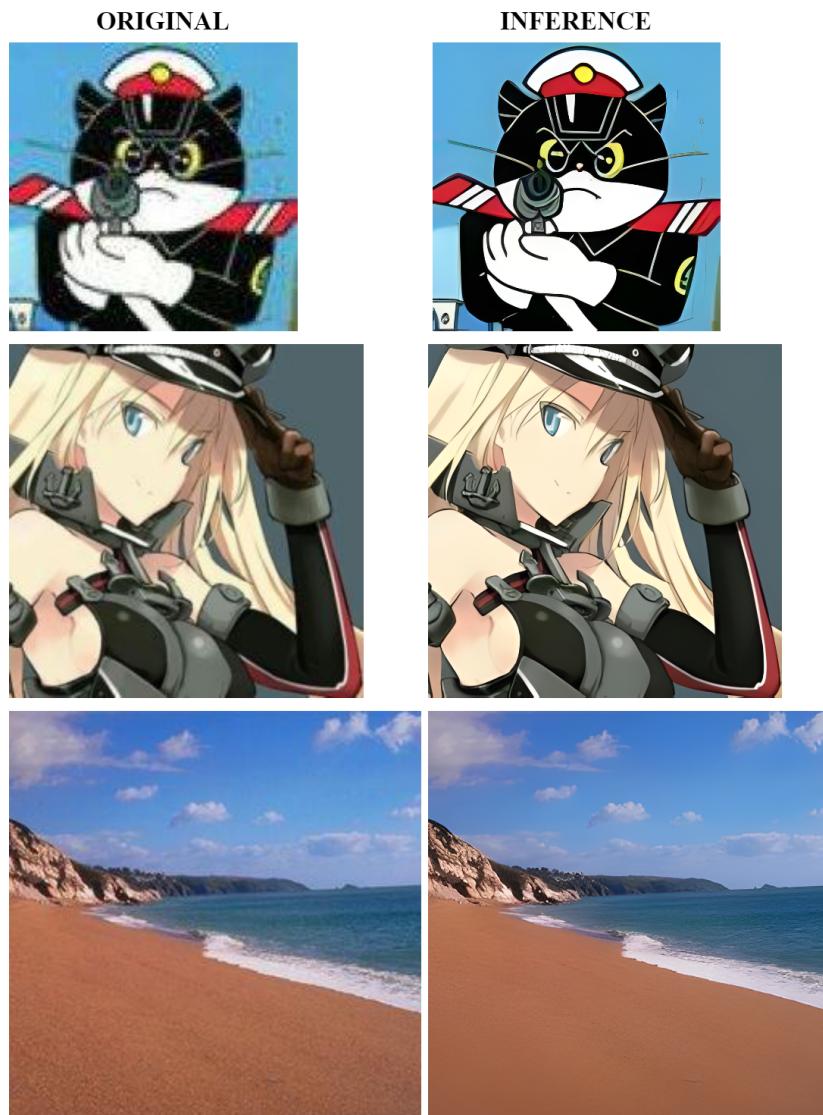


Figure 5. Inference Model

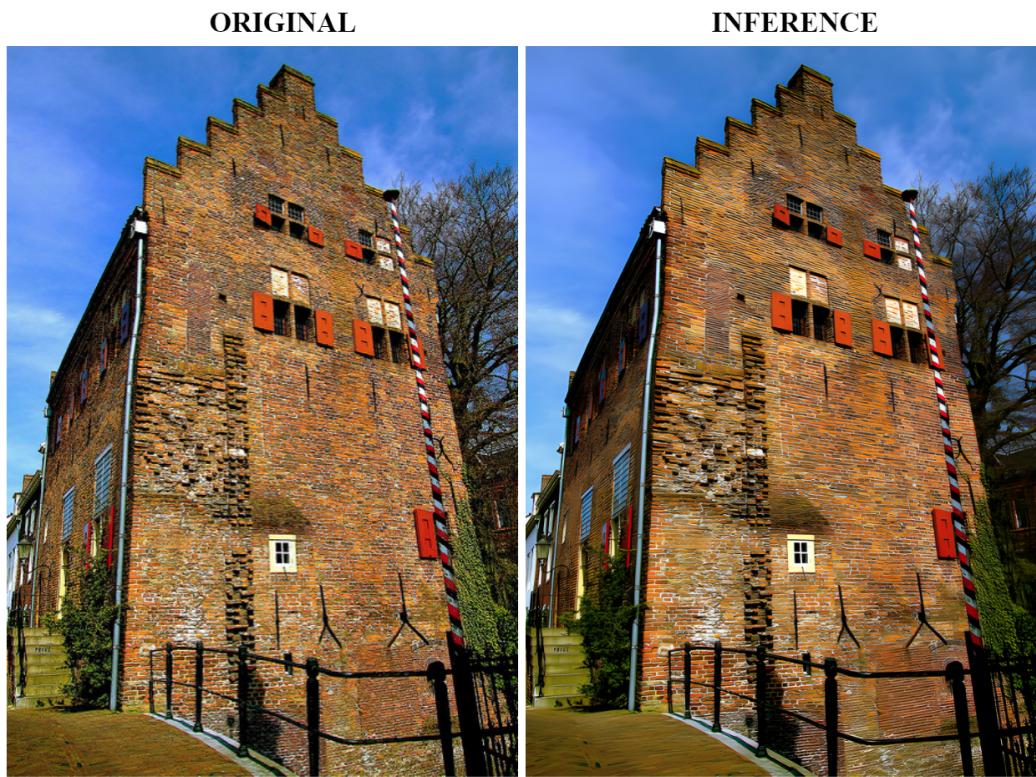


Figure 6. Model Inference



Figure 7. Inference Model



## **9. Conclusion**

Overall, this project was an implementation of the paper “Real-ESRGAN: Training Real-World Blind Super-Resolution with Pure Synthetic Data” [1]. Through this project, we explained and coded a second-order degradation process to model practical degradations and utilize sinc filters to model common ringing and overshoot artifacts. The classical first-order degradation process is not sufficient for the state-of-the-art super-resolution models. U-Net discriminator with spectral normalization to increase discriminator capability and stabilize the training dynamics was also introduced and coded. Finally, Real-ESRGAN was trained from scratch with pure synthetic data (images generated/created and not from real-world images) coming from the three datasets of DIV2K, Flickr2K, and OutdoorSceneTraining (OST). ESRGAN is a combination of three concepts. “Enhanced” refers to the improvement in the quality of images. “Super-resolution” is a technique used to increase the resolution of an image. “Generative Adversarial Network” (GAN) is a type of AI model used for generating new data instances. The results of the Real-ESRGAN are very impressive and hard to believe. I am glad that I was able to reproduce the results and be able to work with cutting-edge technologies and have a model that I can use to enhance my images. In the future, I will retrain the Real-ESRGAN model for the x1 and x2 scale factors. Moreover, I would love to see how higher-order degradation processes such as third-class will affect the results.

## References

- [1] Xintao Wang, Liangbin Xie, Chao Dong, and Ying Shan. Real-ESRGAN: Training Real-World Blind Super-Resolution with Pure Synthetic Data. In ECCV, 2021.
- [2] Chao Dong, Chen Change Loy, and Xiaou Tang. Accelerating the Super-Resolution Convolutional Neural Network. In ECCV, 2016.
- [3] Dong, C., Loy, C.C., He, K., Tang, X.: Learning a deep convolutional network for image super-resolution. In: ECCV, 2014.
- [4] Dong, C., Loy, C.C., He, K., Tang, X.: Image super-resolution using deep convolutional networks. In TPAMI, 2015.
- [5] Wang, Z., Liu, D., Yang, J., Han, W., Huang, T.: Deeply improved sparse coding for image super-resolution. In ICCV, 2015.
- [6] Zihao Wei, Yidong Huang, Yuang Chen, Chenhao Zheng, and Jinnan Gao. A-ESRGAN: Training Real-World Blind Super-Resolution with Attention U-Net Discriminators. 2021.
- [7] Xintao Wang, Ke Yu, Shixiang Wu, Jinjin Gu, Yihao Liu, Chao Dong, Yu Qiao, and Chen Change Loy. Esrgan: Enhanced super-resolution generative adversarial networks. In ECCV, 2018.