

# Technical Appendix

## Algorithms and Parameters

We present the details of the preprocessing functions used in the GYM. Each method is modified or upgrade based on proposed methods with the aid of randomness. The hyperparameters we adopted in the preprocessing procedures for both the intensive preprocessing and the lightweight preprocessing are listed in Table 1.

Notation	Meaning	Value
$\delta$	distortion limit of the Optical Distortion	0.5
$\gamma_1$	Gamma Compression's gamma value	0.6
$\gamma_2$	Gamma Extension's gamma value	2.6
$\sigma$	scale limit of the RSDP	1.3
$T$	translation limit of the SAT	0.16
$S$	sacaling limit of the SAT	0.16
$R$	rotation limit of the SAT	4

Table 1: Hyperparameters' settings during the Preprocessing.

### Optical Distortion

Different from (Liu, Malcolm, and Xu 2010), the Optical Distortion we upgraded and utilized in the GYM is based on assigning a random distortion value chosen from a uniform distribution of the distortion limit. This random process can distort each sample on a different scale for a different time, thus better help the infected model better adapt to the remapping distortions. The details of the Random Pincushion Distortion we proposed and improved in the GYM are explained in Algorithm 1. The random pincushion distortion can be interpreted into three phases. For starters, we acquire a random distortion value,  $\delta_k$ , from a uniform distribution between  $-\delta$  to 0. Using this randomly sampled  $\delta_k$ , we can acquire two pincushion maps for horizontal and vertical indexes, respectively. Finally, by broadcasting those two maps for each pixel, we can output the result. During the experiment, we set the  $\delta$  as 0.5 based on experimental analysis.

### Gamma Compression and Extension

Inspired by the previous work (Kumari, Thomas, and Sahoo 2014), the Gamma Compression and the Gamma Extension are fine-tuned and used in the median filters set to merging pixels' values and enhance the effects of the median filters. The Gamma value of the Gamma Compression procedure is

---

### ALGORITHM 1: Random Pincushion Distortion

---

```

Input: original image  $I \in \mathbb{R}^{h \times w}$ 
Output: distorted image  $I' \in \mathbb{R}^{h \times w}$ 
Parameters: distortion limit  $\delta$ ;

/* 1.Acquire distortion parameter  $\delta_k$  */
1  $\delta_k \sim \mathcal{U}(-\delta, 0)$ ;
/* 2.Acquire Distortion Maps */
2  $c_x = \lfloor (w/2) \rfloor, c_y = \lfloor (h/2) \rfloor$ ;
3  $P_{set} = \{(m, n) \in \{(0, \dots, w) \times (0, \dots, h)\}\}$ ;
4 for  $(u, v)$  in  $P_{set} \setminus \{(m, n)\}$  do
5 |    $map_x(u, v) = ((u - c_x) \times (1 + k)) + c_x$ ;
6 |    $map_y(u, v) = ((v - c_y) \times (1 + k)) + c_y$ ;
7 end
/* 3.Remapping  $I$  to  $I'$  */
8 for  $(u, v)$  in  $P_{set} \setminus \{(m, n)\}$  do
9 |    $I'(u, v) = I(map_x(u, v), map_y(u, v))$ ;
10 end
11 return  $I'$ ;

```

---

set to 0.6, which acquires a Look-Up Table shown in the middle of Figure 1. As demonstrated that larger values from the original pixels range (the left part of 1) are mapping with a larger value close to the maximum value (255), thus helps larger values to bend in. As a result, the median filter can work more efficiently to smoothen pixels of low value. Vice versa, with a Gamma value set to 2.6, we can use the help of the Gamma Extension to merge small values, thus better smoothen large pixels. We summarize the Gamma Compression and Extension as a single function shown in Algorithm 2. As demonstrated, the Gamma Transformation we used here in the experiment can be interpreted as two functional parts. First, we acquire the LUT based on the Gamma value,  $\gamma$ . Then, the output image can be obtained by using the value of the corresponding position in the LUT to replace the original pixel value. The function with a Gamma value larger than 1 conducts extension, and a Gamma value smaller than 1 performs compression. We chose 0.6 and 2.6 as the Gamma values for the compression and extension based on experimental results, as they can achieve better results during the inference after GYM fine-tuning.

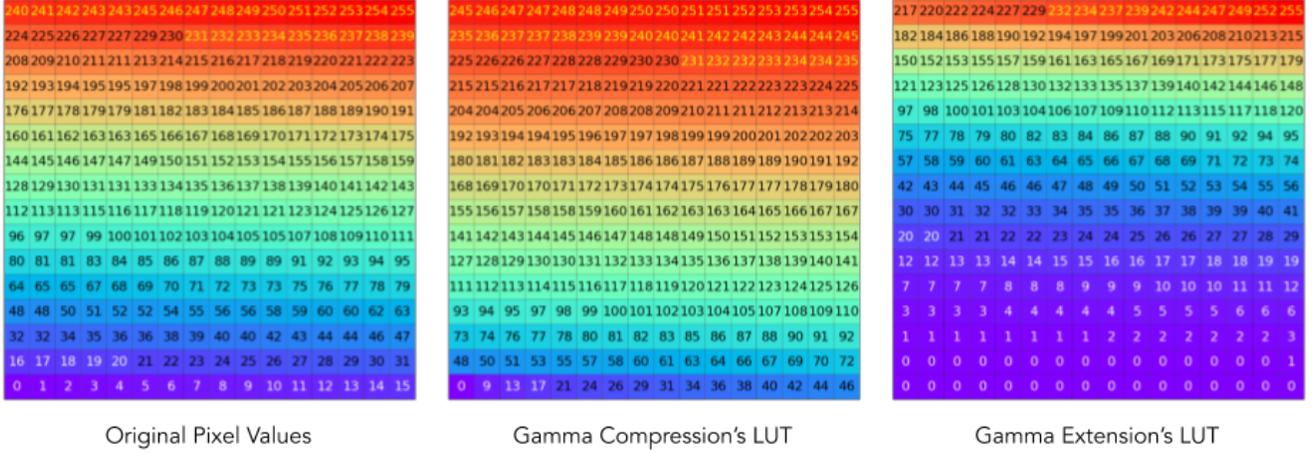


Figure 1: Different Gamma Look-Up Tables (LUTs) used in the Median Filters set: The left is the original pixel values, ranging from 0 to 255; the Gamma Compression uses a Gamma value of 0.6, which lead to the LUT shown in the middle; the Gamma Extension uses a Gamma value of 2.6, which lead to the left LUT.

---

### ALGORITHM 2: Gamma Transformation

---

**Input:** original image  $I \in \mathbb{R}^{h \times w}$   
**Output:** transformed image  $I' \in \mathbb{R}^{h \times w}$   
**Parameters:** Gamma Value  $\gamma$ ;

```

/* 1.Acquire LUT */
1  $T = \text{range}(0 : 255)^{16 \times 16}$ ;
2  $LUT = (T/255)^\gamma \times 255$ ;
/* 2.Assigning New Values */
3  $P_{set} = \{(m, n) \in \{(0, \dots, w) \times (0, \dots, h)\}\}$ ;
4 for  $(u, v)$  in  $P_{set} \setminus \{(m, n)\}$  do
5    $(x, y) = \text{where}(T == I(u, v))$ ;
6    $I'(u, v) = LUT(x, y)$ ;
7 end
8 return  $I'$ ;

```

---



---

### ALGORITHM 3: RSDP

---

**Input:** original image  $I \in \mathbb{R}^{l \times l}$   
**Output:** distorted image  $I' \in \mathbb{R}^{l \times l}$   
**Parameters:** scale limit  $\sigma$ ;

```

/* 1.Acquire random parameter */
1  $Len_{max} = \lfloor (l \times \sigma) \rfloor$ ;
2  $Len \sim \mathcal{U}(l, Len_{max})$ ;
3  $l_{rem} = Len_{max} - Len$ ;
4  $x_1 \sim \mathcal{U}(0, l_{rem})$ ,  $y_1 \sim \mathcal{U}(0, l_{rem})$ ;
5  $x_2 = l_{rem} - x_1$ ,  $y_2 = l_{rem} - y_1$ ;
/* 2.Padding to  $Len_{max} \times l$  */
6  $I' = \text{reshape}(I)$  s.t.  $I' \in \mathbb{R}^{Len \times Len}$ ;
7  $I' = \text{pad}(I', ((x_1, x_2), (y_1, y_2)), \text{value} = 0)$ 
   s.t.  $I' \in \mathbb{R}^{Len_{max} \times Len_{max}}$ ;
/* 3.Reshape  $I'$  to the size of  $I$  */
8  $I' = \text{reshape}(I')$  s.t.  $I' \in \mathbb{R}^{l \times l}$ ;
9 return  $I'$ ;

```

---

## 52 Random Scale Down with Padding (RSDP)

53 We proposed Random Scale Down with Padding (RSDP) as  
54 a tool to help the infected model to better adapt to affine  
55 transformations. The details of the proposed preprocessing  
56 function are explained in 3. The  $\sigma$  we used in the experi-  
57 ment is set to 1.3 to downscale the input image in a range  
58 of range (0.8,1). The whole process of the proposed RSDP  
59 can be interpreted as three functional parts. First, the algo-  
60 rithm acquires random parameters for the scaling and the  
61 padding. This includes after-padding size,  $Len_{max}$ ; resizing  
62 size,  $Len$ ; the number of pixels to pad to reach the after-  
63 padding size,  $l_{rem}$ ; and padding coordinates,  $(x_1, x_2)$  and  
64  $(y_1, y_2)$ . Padding the resized image using the padding coor-  
65 dinates to  $(Len_{max}, Len_{max})$ , we can acquire a black canvas  
66 patched with the resized original input. With resizing the im-  
67 age back to the original size, we can acquire the final result.  
68 Based on the experiment, we found that resizing the image  
69 from 0.8 to 1 times smaller can best help the infected model  
70 adapts to the affine transformation.

## Stochastic Affine Transformation

We adopt the Stochastic Affine Transformation (SAT) (Zeng et al. 2020) to further help the infected model adapt the affine transformation in GYM Fine-tuning and invalidate potential triggers during inference. The details of the SAT are explained as follows in the Algorithm 4. We adopt the same settings from (Zeng et al. 2020),  $T$ , 0.16,  $S$ , 0.16, and  $R$ , 4. The whole process of the SAT can be functionally interpreted as three parts, namely translation, rotation, and scaling.

## Attacking and Preprocessing Results

The attacking and preprocessing results over patched sam-  
ples are shown in Figure 2,3,and 4. Please noted that the  
intensive preprocessed patched data are not used during the  
GYM fine-tuning procedure. We only adopt the intensive pre-  
processed patched data during the inference before and after  
the fine-tuning (Inference(I)).

71  
72  
73  
74  
75  
76  
77  
78  
79  
80  
81  
82  
83  
84  
85  
86

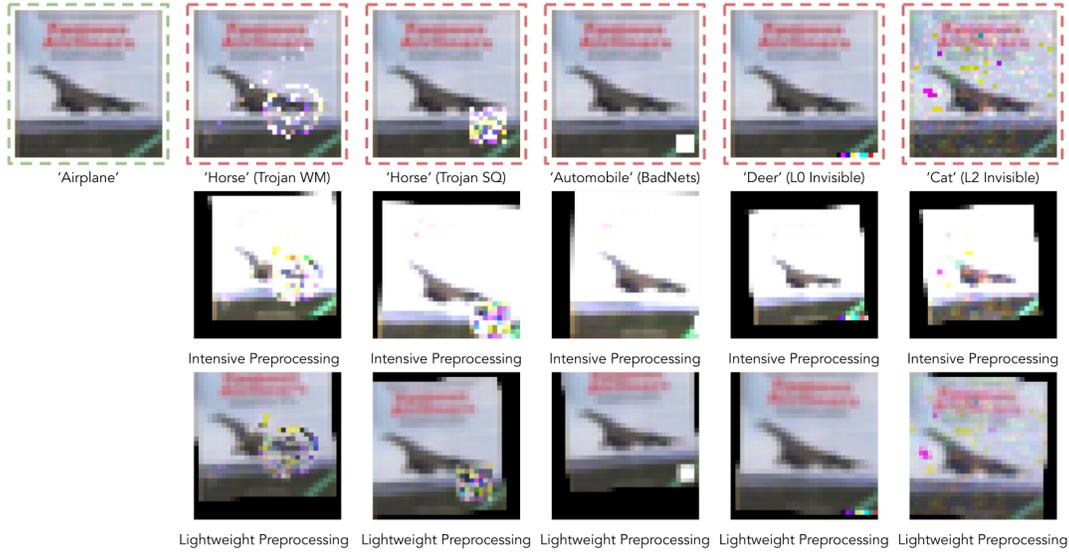


Figure 2: The attacking and preprocessing results over the Cifar10 dataset.

#### ALGORITHM 4: SAT

---

**Input:** original image  $I \in \mathbb{R}^{h \times w}$   
**Output:** transformed image  $I' \in \mathbb{R}^{h \times w}$   
**Parameters:** translation limit  $T$ ; scaling limit  $S$ , rotation limit  $R$ .

```

1  $I' = O^{h \times w}$ ;
  /* 1.Translation */
2  $\delta_x \sim \mathcal{U}(-T, T)$ ;
3  $\delta_y \sim \mathcal{U}(-T, T)$ ;
4  $\Delta_x = \delta_x \times w$ ;
5  $\Delta_y = \delta_y \times h$ ;
6 if  $(x + \Delta_x \in (0, w)) \wedge (y + \Delta_y \in (0, h))$  then
7    $I'(x, y) = I(x + \Delta_x, y + \Delta_y)$ ;
8 end
  /* 2.Rotation */
9  $\delta_r \sim \mathcal{U}(-R, R)$ ;
10  $\Delta_r = \delta_r \times \pi/180$ ;
11 for  $(x_i, y_j)$  in  $\{(x, y) | x \in (0, w), y \in (0, h)\}$  do
12    $x'_i = -(x_i - \lfloor w/2 \rfloor) \times \sin(\Delta_r) + (y_j - \lfloor h/2 \rfloor) \times \cos(\Delta_r)$ ;
13    $y'_j = (x_i - \lfloor w/2 \rfloor) \times \cos(\Delta_r) + (y_j - \lfloor h/2 \rfloor) \times \sin(\Delta_r)$ ;
14    $x'_i = \lfloor x'_i + \lfloor w/2 \rfloor \rfloor$ ;
15    $y'_j = \lfloor y'_j + \lfloor h/2 \rfloor \rfloor$ ;
16   if  $(x'_i \in (0, w)) \wedge (y'_j \in (0, h))$  then
17      $I'(x_i, y_j) = I(x'_i, y'_j)$ ;
18   end
19 end
  /* 3.Scaling */
20  $\delta_s \sim \mathcal{U}(1 - S, 1 + S)$ ;
21  $h_{new} = \delta_s \times h$ ;
22  $w_{new} = \delta_s \times w$ ;
23  $I' = \text{reshape}(I', (h_{new}, w_{new}))$ ;
24 if  $\delta_s > 1$  then
25    $I'(x, y) = \text{cropping}(I', (h, w))$ ;
26 end
27 if  $\delta_s < 1$  then
28    $I'(x, y) = \text{padding}(I', (h, w))$ ;
29 end
30 return  $I'$ ;

```

---

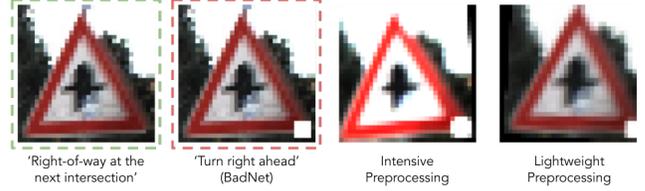


Figure 3: The attacking and preprocessing results over the GTSRB dataset.

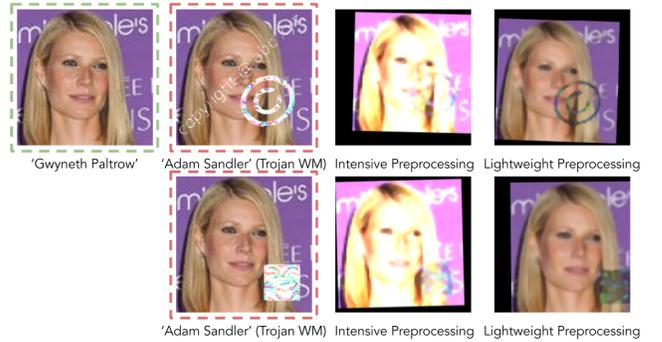


Figure 4: The attacking and preprocessing results over the PubFig dataset.

As shown in each figure, the triggers in each sample pre-processed by the intensive preprocess become hard to recognize by human eyes. Thus, it is intuitive that such intensive preprocessing can help the infected models revise their decision boundaries to encompass those patch data to their original classes. With the decision boundary being shifted with the intensive preprocessed data, the lightweight preprocessing can help the fine-tuned model achieve a more accurate result, with less preprocessing procedures being adopted.

87  
88  
89  
90  
91  
92  
93  
94  
95

96 Please noted that the lightweight preprocessing is actually  
97 proposed based on the intensive preprocessing, which sim-  
98 plified some steps in the intensive preprocessing. The inten-  
99 sive preprocessing is developed from a functional perspec-  
100 tive to help the infected model better and faster adapt to the  
101 lightweight preprocessing, thus achieving high accuracy for  
102 an efficient deployment.

### 103 **References**

- 104 Kumari, A.; Thomas, P. J.; and Sahoo, S. 2014. Single image  
105 fog removal using gamma transformation and median filter-  
106 ing. In *2014 annual IEEE India conference (INDICON)*, 1–5.  
107 IEEE.
- 108 Liu, T.; Malcolm, A.; and Xu, J. 2010. Pincushion distor-  
109 tion correction in x-ray imaging with an image intensifier. In  
110 *Fourth International Conference on Experimental Mechan-*  
111 *ics*, volume 7522, 75223T. International Society for Optics  
112 and Photonics.
- 113 Zeng, Y.; Qiu, H.; Memmi, G.; and Qiu, M. 2020. A  
114 Data Augmentation-based Defense Method Against Ad-  
115 versarial Attacks in Neural Networks. *arXiv preprint*  
116 *arXiv:2007.15290* .