

# Appendix: Out-of-Distribution Detection with Semantic Mismatch under Masking

## A Training Process of Generative Model

### A.1 Objective Functions

We implement the adversarial loss with a U-net based discriminator [4], denoted as  $D^{Unet}$ .  $D^{Unet}$  contains two components:  $D_{enc}^{Unet}$  and  $D_{dec}^{Unet}$ .  $D_{enc}^{Unet} \in \mathbb{R}$  provides the real/fake decision as a scalar. While  $D_{dec}^{Unet} \in \mathbb{R}^I$  generates a per-pixel real/fake map for the input image, where  $I = h \times w$  indicates the scale of input image. Compared with the vanilla discriminator,  $D^{Unet}$  not only determines whether the input image is realistic or fake, but also tries to locate the fake parts. Empowered by the per-pixel real/fake map, our generative model can be optimized to focus more on structural semantic features and synthesize coherent image both globally and locally as desired. We formulate the adversarial loss for the discriminator in Eq. (1)-Eq. (3):

$$\mathcal{L}_{D^{Unet}} = \mathcal{L}_{D_{enc}^{Unet}} + \mathcal{L}_{D_{dec}^{Unet}}, \quad (1)$$

$$\mathcal{L}_{D_{enc}^{Unet}} = -\mathbb{E}_x[\log D_{enc}^{Unet}(x, y)] - \mathbb{E}_x[\log(1 - D_{enc}^{Unet}(x', y))], \quad (2)$$

$$\mathcal{L}_{D_{dec}^{Unet}} = -\mathbb{E}_x\left[\sum_I \log D_{dec}^{Unet}(x, y)\right] - \mathbb{E}_x\left[\sum_I \log(1 - D_{dec}^{Unet}(x', y))\right], \quad (3)$$

where  $\mathcal{L}_{D^{Unet}}$  and  $\mathcal{L}_{D_{dec}^{Unet}}$  are the loss functions for  $D_{enc}^{Unet}$  and  $D_{dec}^{Unet}$ , respectively. Correspondingly, the adversarial loss applied on the generator is as follow:

$$\mathcal{L}_G = -\mathbb{E}_x\left[\log D_{enc}^{Unet}(x', y) + \sum_I \log D_{dec}^{Unet}(x')\right] + \ell_1(x, x') + \ell_2(x, x') + \mathcal{SSIM}(x, x'). \quad (4)$$

### A.2 Training Process

**Encoder.** We adopt a four-layer convolutional neural network as the feature extractor for Encoder, then two fully-connected layers are employed to output  $\mu$  and  $\Sigma$ . The dimension of latent variable  $z$  is set at 128.

**Decoder.** We employ the generator architecture proposed in [1] as our Decoder's backbone, then reset the input size to (3, 32, 32), and the channel multiplier to 32, which represents the number of units in each layer [1]. The input latent variable size equals 128.

**Discriminator.** We build up  $D^{Unet}$  based on the implementation of [4], changing the channel multiplier to 32.

All three models mentioned above are trained from scratches in an end-to-end way. We use Adam [3] as the optimizer, with  $\beta_1 = 0$ ,  $\beta_2 = 0.999$ , learning rate fixed at  $5 \cdot 10^{-5}$ . The batch size is set at 96. We detail the training process of our generative model in Algorithm 1.

---

**Algorithm 1: Training Framework of G**

---

**Input** : Training data  $\mathcal{X} = \{x\}^N$ ,  $\mathcal{Y} = \{y\}^N$ , the random mask  $\mathbf{M}$   
**Output** : The parameters of  $\mathbf{E}$ ,  $\mathbf{D}$

- 1 **for** *some training iterations* **do**
- 2      $x' = \mathbf{G}(\mathbf{M}(x), y) = \mathbf{D}(\mathbf{E}(\mathbf{M}(x), y))$ ;
- 3     Feed  $(x, y)$  and  $(x', y)$  into  $D^{Unet}$ , **respectively**;
- 4     Optimize  $\mathbf{D}$  and  $\mathbf{E}$  for  $\mathcal{L}_{\mathbf{G}}$  (Eq. (4)) and  $\mathcal{L}_{KLD}$ ;
- 5     Optimize  $D^{Unet}$  for  $\mathcal{L}_{D^{Unet}}$  (Eq. (1));
- 6 **end**
- 7 **return**  $\mathbf{E}$ ,  $\mathbf{G}$

---

## B Quantitative Results

In this section, we provide more experimental results on CIFAR-10 and CIFAR-100 benchmarks, respectively. In addition, to further validate the effectiveness of the proposed *conditional binary classifier* ( $\mathcal{C}_b$ ) in anomalous scoring model, we detail its performance on each OOD dataset by varying the type of  $\mathcal{C}_b$ , i.e. trained with/without external OOD data.

### B.1 More Results on CIFAR-10 Benchmarks

Table 1 presents the comparison of our MOODCAT trained with external unlabeled data sourcing from TINY-IMAGENET, and baselines implemented with extra data. We conclude that MOODCAT outperforms or at least on par with baselines on CIFAR-10 benchmarks.

Additionally, in Table 1 we observe that OE and UDG achieve a much better performance on SVHN than on other OOD datasets. In fact, most street number images contained in SVHN have relatively flat backgrounds, as shown in Fig. 3 and Fig. 4’s SVHN columns. In this case, OE and UDG can achieve excellent performance by overfitting to this specific low-level feature of SVHN instead of considering the semantic level shift caused by SVHN. Thus, when encounter a more challenging case, e.g., CIFAR-100, which has the same data source as CIFAR-10 but different semantic meanings, both OE and UDG suffer a noticeable performance degradation. In contrast, MOODCAT identifies OOD according to their semantic mismatch, thus, remains stable performance on various OODs.

Table 1: OOD Detection Performance on CIFAR-10 benchmarks, MOODCAT trained with external OOD data. All the values are in percentages.  $\uparrow/\downarrow$  indicates higher/lower value is better. The best results are in **bold**.

Detection Methods	OOD	FPR@ TPR95% $\downarrow$	AUROC $\uparrow$	AUPR In $\uparrow$	AUPR Out $\uparrow$	Classification Accuracy $\uparrow$
MCD	SVHN	60.27	89.78	85.33	94.25	90.56
	CIFAR-100	74.00	82.78	83.97	79.16	90.56
	TINY-IMAGENET	78.89	80.98	85.63	72.48	87.33
	TEXTURE	83.92	81.59	90.20	63.27	90.56
	LSUN	68.96	84.71	85.74	81.50	90.56
	PLACES365	72.08	83.51	69.44	92.52	88.51
	<b>Mean</b>	73.02	83.89	83.39	80.53	89.68
OE	SVHN	20.88	96.43	93.62	98.32	91.87
	CIFAR-100	58.54	86.22	86.17	84.88	91.87
	TINY-IMAGENET	58.98	87.65	90.09	82.16	89.27
	TEXTURE	51.17	89.56	93.79	81.88	91.87
	LSUN	57.97	86.75	87.69	85.07	91.87
	PLACES365	55.64	87.00	73.11	94.67	90.99
	<b>Mean</b>	50.53	88.93	87.55	87.83	91.29
UDG	SVHN	<b>13.26</b>	<b>97.49</b>	<b>95.66</b>	<b>98.69</b>	92.94
	CIFAR-100	47.20	90.98	<b>91.74</b>	89.36	92.94
	TINY-IMAGENET	50.18	91.91	<b>94.43</b>	86.99	90.22
	TEXTURE	20.43	96.44	98.12	92.91	92.94
	LSUN	42.05	93.21	94.53	91.03	92.94
	PLACES365	44.22	92.64	87.17	96.66	91.68
	<b>Mean</b>	36.22	93.78	93.61	92.61	92.28
Ours	SVHN	24.27	95.93	92.98	98.05	<b>95.02</b>
	CIFAR-100	<b>39.92</b>	<b>91.45</b>	91.54	<b>91.73</b>	<b>95.02</b>
	TINY-IMAGENET	<b>32.41</b>	<b>93.34</b>	93.63	<b>93.41</b>	<b>92.54</b>
	TEXTURE	<b>6.86</b>	<b>98.69</b>	<b>99.29</b>	<b>97.71</b>	<b>95.02</b>
	LSUN	<b>33.31</b>	<b>93.40</b>	<b>93.85</b>	<b>93.22</b>	<b>95.02</b>
	PLACES365	<b>35.51</b>	<b>92.77</b>	<b>82.25</b>	<b>94.82</b>	<b>93.87</b>
	<b>Mean</b>	<b>28.71</b>	<b>94.27</b>	<b>92.26</b>	<b>94.82</b>	<b>94.42</b>

## B.2 More Results on CIFAR-100 benchmarks

Table 2 shows the comparison of our MOODCAT trained without external OOD data, and baselines are implemented under the same setting. We conclude that MOODCAT achieve the state-of-the-art performance on CIFAR-100 benchmarks.

## B.3 Ablation Study on Conditional Binary Classifier

To study how much the proposed *Conditional Binary Classifier* ( $\mathcal{C}_b$ ) contributes to MOODCAT, we conduct several ablations on  $\mathcal{C}_b$ . More specific, we consider three configurations:  $\mathcal{C}_b$ ,  $\mathcal{C}_b(\text{TINY-IMAGENET})$ , and  $\mathcal{C}_b + \mathcal{C}_b(\text{TINY-IMAGENET})$ , where  $\mathcal{C}_b$  referring to the Conditional Binary Classifier trained only on In-D samples,  $\mathcal{C}_b(\text{TINY-IMAGENET})$  denoted the Conditional Binary Classifier using TINY-IMAGENET as extra training data and  $\mathcal{C}_b + \mathcal{C}_b(\text{TINY-IMAGENET})$  indicating that  $\mathcal{C}_b$  and  $\mathcal{C}_b(\text{TINY-IMAGENET})$  are used in a cascade way.

Table 3 and Table 4 demonstrate  $\mathcal{C}_b$ 's performance on CIFAR-10 and CIFAR-100 benchmarks across six OOD datasets, respectively. The main takeaways

Table 2: OOD Detection Performance on CIFAR-100 as In-D, MOODCAT training without external data. All the values are in percentages.  $\uparrow/\downarrow$  indicates higher/lower value is better. The best results are in **bold**.

Detection Methods	OOD	FPR@ TPR95% $\downarrow$	AUROC $\uparrow$	AUPR In $\uparrow$	AUPR Out $\uparrow$	Classification Accuracy $\uparrow$
ODIN	SVHN	90.33	75.59	65.25	84.49	76.65
	CIFAR-10	81.28	77.90	79.93	73.39	76.65
	TINY-IMAGENET	82.74	77.58	86.26	61.38	69.56
	TEXTURE	79.47	77.92	86.69	62.97	76.65
	LSUN	80.57	78.22	86.34	63.44	76.10
	PLACES365	76.42	80.66	66.77	89.66	77.56
	<b>Mean</b>	81.89	77.98	78.54	72.56	75.53
EBO	SVHN	78.23	83.57	75.61	90.24	76.65
	CIFAR-10	81.25	78.95	80.01	74.44	76.65
	TINY-IMAGENET	83.32	78.34	87.08	62.13	69.56
	TEXTURE	84.29	76.32	85.87	59.12	76.65
	LSUN	84.51	77.66	86.42	61.40	76.10
	PLACES365	78.37	80.99	68.22	89.60	77.56
	<b>Mean</b>	81.66	79.31	80.54	72.82	75.53
Ours	SVHN	<b>58.16</b>	<b>87.38</b>	<b>78.25</b>	<b>93.81</b>	<b>76.65</b>
	CIFAR-10	<b>54.31</b>	<b>85.91</b>	<b>86.27</b>	<b>85.91</b>	<b>76.65</b>
	TINY-IMAGENET	<b>55.33</b>	<b>86.95</b>	<b>87.55</b>	<b>86.67</b>	<b>69.56</b>
	TEXTURE	<b>46.70</b>	<b>89.20</b>	<b>93.48</b>	<b>83.28</b>	<b>76.65</b>
	LSUN	<b>53.43</b>	<b>87.98</b>	<b>88.82</b>	<b>87.32</b>	<b>76.10</b>
	PLACES365	<b>54.20</b>	<b>87.41</b>	<b>71.68</b>	<b>95.78</b>	<b>77.56</b>
	<b>Mean</b>	<b>53.69</b>	<b>87.47</b>	<b>84.34</b>	<b>88.80</b>	<b>75.53</b>

are: (1)  $C_b$  or  $C_b(\text{TINY-IMAGENET})$  alone can achieve acceptable performance; (2)  $C_b(\text{TINY-IMAGENET})$  outperforms  $C_b$ , which means that adding external unlabeled data into the training process can improve the detection ability; (3) coupling scorers, here  $C_b + C_b(\text{TINY-IMAGENET})$ , usually leads to a better detection capability than that of any single scorer within the coupling. Above findings align with what we have reported in our paper, and further indicate that  $C_b$  plays a key role in the proposed anomalous scoring model.

#### B.4 Ablation Study on Masking Style

We try several masking forms as exemplified in Fig. 1, and summarize corresponding experimental results in Table 5. Experiments show the randomly masking outperforms other strategies.

From the first three rows in Table 5, we notice that masking can indeed help with performance improvement. However, as we can observed from the second column in Fig. 1, a fixed mask with high ratio (e.g., 0.3) can lead the synthesis to loss fine details. In addition, we implement a patched masking like [2]. However, such masking style may break the continuity within the image, thus lead to low quality on the synthesis for In-D. We also try a non-masking strategy, shuffling, but it further break the continuity of the image. Finally, we identify that the most effective strategy is randomly masking. As can be seen, both the quality for synthesis in Fig. 1 and overall performance in Table 5 outperform other strategies.

Table 3: Conditional Binary Classifier Performance on CIFAR-10 benchmarks. All the values are in percentages.  $\uparrow/\downarrow$  indicates higher/lower value is better. The best results are in **bold**.  $\mathcal{C}_b$  and  $\mathcal{C}_b(\text{TINY-IMAGENET})$  indicates the proposed model trained without/with external unlabeled TINY-IMAGENET data, respectively.

Anomalous Scoring Model	OOD	FPR@ TPR95% $\downarrow$	AUROC $\uparrow$	AUPR In $\uparrow$	AUPR Out $\uparrow$
$\mathcal{C}_b$	SVHN	48.01	86.85	75.20	94.34
	CIFAR-100	42.80	89.13	88.58	89.85
	TINY-IMAGENET	40.54	89.78	89.27	90.42
	TEXTURE	42.54	87.47	91.33	83.85
	LSUN	43.76	90.15	90.18	90.17
	PLACES365	43.49	89.40	72.82	96.65
	<b>Mean</b>	43.52	88.80	84.56	90.88
$\mathcal{C}_b(\text{TINY-IMAGENET})$	SVHN	39.47	91.49	83.58	96.23
	CIFAR-100	37.43	91.31	91.02	91.73
	TINY-IMAGENET	31.92	93.10	93.01	93.34
	TEXTURE	25.74	94.25	96.17	91.89
	LSUN	32.74	93.55	93.83	93.40
	PLACES365	34.45	92.78	81.48	97.71
	<b>Mean</b>	33.63	92.75	89.85	94.05
$\mathcal{C}_b + \mathcal{C}_b(\text{TINY-IMAGENET})$	SVHN	<b>39.44</b>	<b>91.50</b>	<b>83.60</b>	<b>96.25</b>
	CIFAR-100	<b>36.64</b>	<b>91.40</b>	<b>91.15</b>	<b>91.85</b>
	TINY-IMAGENET	<b>31.86</b>	<b>93.12</b>	<b>93.04</b>	<b>93.38</b>
	TEXTURE	<b>25.37</b>	<b>94.34</b>	<b>96.24</b>	<b>92.00</b>
	LSUN	<b>32.67</b>	<b>93.55</b>	<b>93.84</b>	<b>93.41</b>
	PLACES365	<b>34.42</b>	<b>92.79</b>	<b>81.51</b>	<b>97.72</b>
	<b>Mean</b>	<b>33.40</b>	<b>92.78</b>	<b>89.90</b>	<b>94.10</b>

## B.5 Experiments on Advanced Classifier architectures

We empower UDG with wider (WRN28) and deeper (DenseNet) classifier. Table 6 shows the comparison results with CIFAR-100 as In-D samples using WRN28 and DenseNet architecture.

As can be observed from the table, while UDG performs better on these architectures when compared to ResNet18, it still lags far behind our results.

## C Qualitative Results

In this section, we demonstrate several batches of visual examples of MOODCAT including both In-D and OOD cases.

***In-D samples with their syntheses.*** Fig. 2 visualizes In-D samples and their corresponding syntheses from CIFAR-10 and CIFAR-100, respectively. Note that we expect the syntheses to resemble the input images for In-D samples with correct labels.

***OOD samples with their syntheses.*** Fig. 3 visualizes OOD samples from six datasets, which are employed in the CIFAR-10 benchmarks, along with their corresponding masked images and the syntheses generated by our MOODCAT.

Table 4: Conditional Binary Classifier Performance on CIFAR-100 benchmarks. All the values are in percentages.  $\uparrow/\downarrow$  indicates higher/lower value is better. The best results are in **bold**.  $\mathcal{C}_b$  and  $\mathcal{C}_b(\text{TINY-IMAGENET})$  indicates the proposed model trained without/with external unlabeled TINY-IMAGENET data, respectively.

Anomalous Scoring Model	OOD	FPR@ TPR95% $\downarrow$	AUROC $\uparrow$	AUPR In $\uparrow$	AUPR Out $\uparrow$
$\mathcal{C}_b$	SVHN	65.18	81.32	65.61	91.35
	CIFAR-10	55.11	85.75	85.78	85.99
	TINY-IMAGENET	54.69	86.27	86.26	86.43
	TEXTURE	56.63	83.30	88.40	77.17
	LSUN	54.77	86.96	87.20	86.83
	PLACES365	54.18	86.36	67.60	95.54
	<b>Mean</b>	56.76	84.99	80.14	87.22
$\mathcal{C}_b(\text{TINY-IMAGENET})$	SVHN	54.61	<b>86.30</b>	<b>74.30</b>	93.80
	CIFAR-10	49.82	87.57	87.74	87.69
	TINY-IMAGENET	45.86	89.38	89.48	89.43
	TEXTURE	48.24	87.16	91.55	81.83
	LSUN	44.43	90.07	90.25	90.00
	PLACES365	46.89	88.93	72.99	96.41
	<b>Mean</b>	48.31	88.24	84.39	89.86
$\mathcal{C}_b + \mathcal{C}_b(\text{TINY-IMAGENET})$	SVHN	<b>54.31</b>	<b>86.30</b>	<b>74.30</b>	<b>93.81</b>
	CIFAR-10	<b>49.62</b>	<b>87.60</b>	<b>87.77</b>	<b>87.77</b>
	TINY-IMAGENET	<b>45.46</b>	<b>89.39</b>	<b>89.48</b>	<b>89.48</b>
	TEXTURE	<b>47.18</b>	<b>87.37</b>	<b>91.71</b>	<b>82.17</b>
	LSUN	<b>44.01</b>	<b>90.08</b>	<b>90.26</b>	<b>90.04</b>
	PLACES365	<b>46.73</b>	<b>88.95</b>	<b>73.02</b>	<b>96.43</b>
	<b>Mean</b>	<b>47.89</b>	<b>88.28</b>	<b>84.42</b>	<b>89.95</b>

In Fig. 4, the In-D dataset changes to CIFAR-100. We employed OOD samples sourced from the same six OOD datasets as that of CIFAR-100 benchmarks in Fig. 4. Note that, when OOD is fed to MOODCAT, we prefer to have a clear distinction between the synthesis generated by MOODCAT and the input image.

## D Further Discussion

### D.1 Computational Cost Analysis

MOODCAT is designed as an auxiliary model, which works in parallel with the classifier. This auxiliary architecture ensures MOODCAT a plug-and-play model without compromising classifier’s accuracy. Meanwhile, MOODCAT can satisfy high performance requirements in the context of OOD detection. However, as an auxiliary model, MOODCAT inevitably introduces extra computation and memory costs.

Table 7 summarizes the computational cost of MOODCAT, and that of ODIN, i.e., ResNet18, and that of widely adopted classifier architectures, ResNet18, WResNet28, WResNet101 in terms of number of multiply-add operations (MAC), and number of model’s parameters (Params). As can be observed, the cost of basic version of MOODCAT, i.e.  $\mathcal{C}_b$ , **E**, **D**, is relatively small, Params 4.552M, MACs

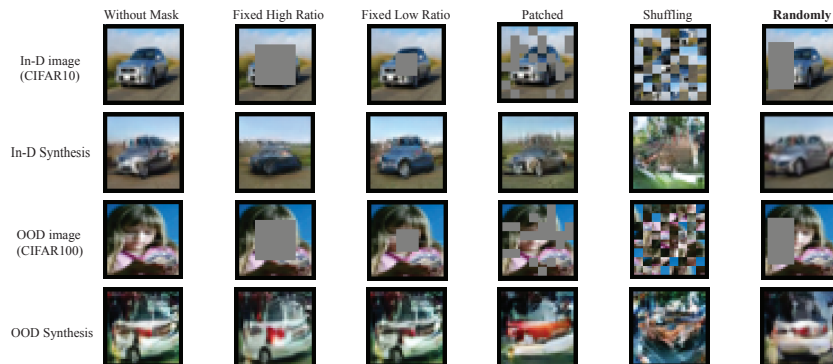


Fig. 1: Visualization of different masking styles and their impacts on synthesized images. The semantic label is assigned as “car” for both the In-D image and the OOD image. We set the masking ratio as 0.3 for “Fixed High Ratio” and “Patched”, 0.1 for “Fixed Low Ratio”, and that of “Randomly” varies from 0.1 to 0.3. MOODCAT employs the **Randomly** masking style.

Table 5: Ablation studies on different masking styles. The results are obtained by setting CIFAR-10 as In-D, CIFAR-100 as OOD, with MOODCAT trained on extra TINY-IMAGENET acting as OOD. The **bolded** values are the highest performance. All the values are in percentages.  $\uparrow/\downarrow$  indicates higher/lower value is better.

Mask Style	FPR@TPR95% $\downarrow$	AUROC $\uparrow$	AUPR-In $\uparrow$	AUPR-Out $\uparrow$
Without Masking	40.53	91.26	91.25	91.55
Fixed Low Ratio	40.20	91.33	91.32	91.59
Fixed High Ratio	39.57	91.56	91.48	91.87
Patched	39.81	91.34	91.33	91.64
Shuffling	44.14	88.73	88.15	89.18
<b>Randomly</b>	<b>39.48</b>	<b>91.66</b>	<b>91.65</b>	<b>91.95</b>

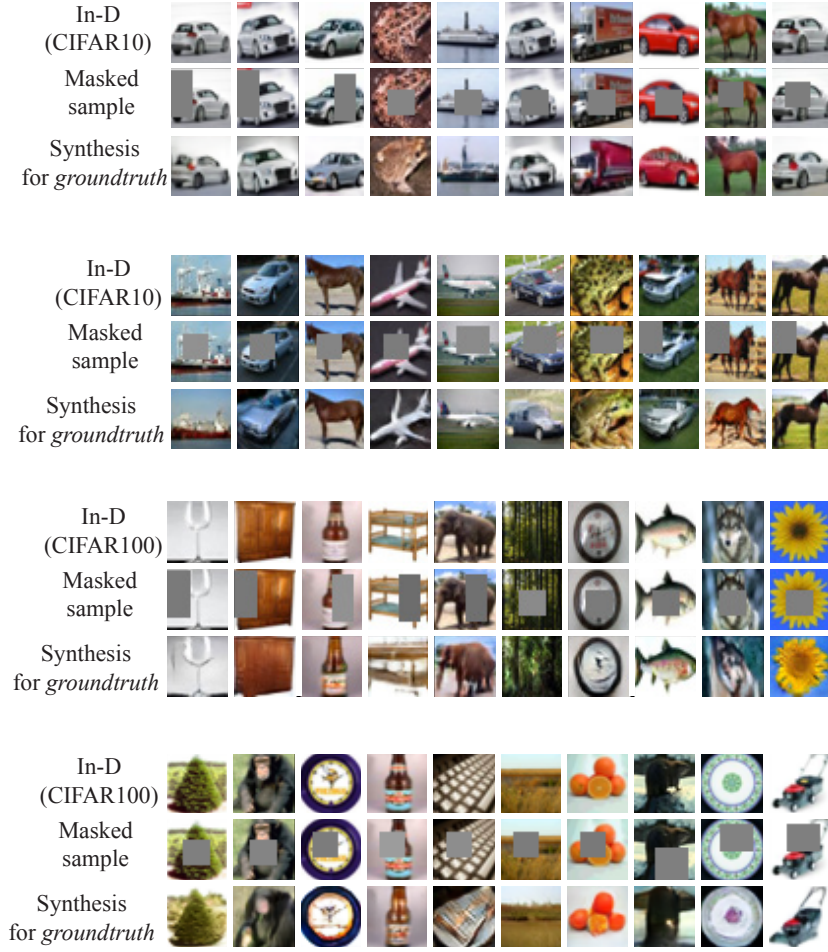


Fig. 2: Visualization results of MOODCAT with CIFAR-10/ CIFAR-100 as In-D. We exemplify several In-D samples in each panel’s first row, following the intermediate masked version, and the last row presents their corresponding synthetic version generated by MOODCAT with the groundtruth labels.



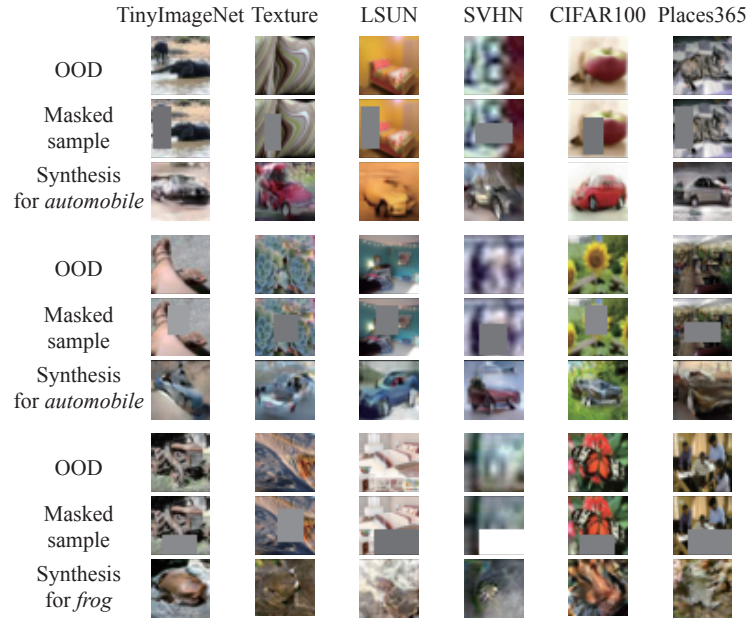


Fig. 3: OOD visualization results of MOODCAT trained on CIFAR-10. In each panel, we exemplify OOD samples across six OOD datasets in the first row, following is the intermediate masked version, the last row presents their corresponding synthetic version generated by MOODCAT with the given semantic label, the same below.

Table 6: Experiments on advanced model architectures. Performance comparison with UDG on CIFAR-100 benchmarks. For our method, we use the results in the main paper with a ResNet18 classifier. We give advantage to UDG, which is reimplemented with deeper/wider WideResNet-28, DenseNet, while MOODCAT’s parameter number is equivalent to ResNet18. **Bold** are the best.

Architecture	OOD dataset	FPR@TPR95 ↓	AUROC ↑	AUPR In ↑	AUPR Out ↑
WideResNet28 UDG	SVHN	66.76	85.29	76.14	92.33
	CIFAR-10	82.35	76.67	78.52	72.63
	TINY-IMAGENET	78.91	79.04	87.00	65.06
	TEXTURE	73.62	79.01	85.53	67.08
	LSUN	77.04	79.79	87.49	66.93
	PLACES365	72.25	81.49	66.72	90.65
	<b>Mean<math>\pm</math>Std</b>	<b>75.16<math>\pm</math>5.49</b>	<b>80.22<math>\pm</math>2.93</b>	<b>80.23<math>\pm</math>8.11</b>	<b>75.78<math>\pm</math>12.44</b>
DenseNet UDG	SVHN	80.67	75.54	75.65	70.99
	CIFAR-10	85.87	74.06	77.16	68.90
	TINY-IMAGENET	82.36	76.81	85.76	61.56
	TEXTURE	76.32	78.93	63.79	89.02
	LSUN	79.12	78.91	66.83	88.23
	PLACES365	73.59	76.27	82.76	65.20
	<b>Mean<math>\pm</math>Std</b>	<b>79.66<math>\pm</math>4.36</b>	<b>76.75<math>\pm</math>1.92</b>	<b>75.33<math>\pm</math>8.64</b>	<b>73.98<math>\pm</math>11.79</b>
MOODCAT (Ours, Res18)	SVHN	<b>51.6</b>	<b>88.99</b>	<b>80.89</b>	<b>94.81</b>
	CIFAR-10	<b>50.17</b>	<b>87.76</b>	<b>88.18</b>	<b>87.79</b>
	TINY-IMAGENET	<b>46.07</b>	<b>89.42</b>	<b>89.73</b>	<b>89.28</b>
	TEXTURE	<b>42.22</b>	<b>90.56</b>	<b>94.43</b>	<b>85.13</b>
	LSUN	<b>47.85</b>	<b>89.96</b>	<b>90.33</b>	<b>89.23</b>
	PLACES365	<b>47.72</b>	<b>89.3</b>	<b>74.83</b>	<b>96.48</b>
	<b>Mean<math>\pm</math>Std</b>	<b>47.61<math>\pm</math>3.29</b>	<b>89.33<math>\pm</math>9.95</b>	<b>86.4<math>\pm</math>7.19</b>	<b>90.45<math>\pm</math>4.33</b>

0.408G, when compared to that of ResNet18 (Params 11.174M, MAC 0.556G) and other widely adopted classifier architectures, e.g., WResNet28 (Params 36.479M, MACs 5.248G). Note that the performance of basic MOODCAT, whose anomalous scoring model only contains  $C_b$ , is still acceptable as shown in Table 3 and Table 4. Thus, if the computational cost is a real concern in the practice, the deployer can adopt MOODCAT with  $C_b$  alone as the anomalous scorer. For the MOODCAT supported by IQA models, e.g., LPIPS, DISTS, the total computational cost is comparable to that of WResNet28 or WResNet101, yet slightly larger than ResNet18. Thus, if the detection ability is put at the first place, one can explore to enhance the anomalous scoring model by employing extra IQA models. Actually, there is a trade-off between the OOD detection performance and computational cost of MOODCAT, and our anomalous scoring model leaves design space for the deployer to explore according to the real-world application.

Table 7: Computational and memory costs of MOODCAT and its components.

Model	E	D	$C_b$	MOODCAT basic	LPIP /DISTS	MOODCAT	ResNet 18	WResNet 28	WResNet 101
Params (M)	0.460	3.821	0.271	<b>4.552</b>	14.715	33.982	11.174	36.479	126.89
MACs (G)	0.0049	0.297	0.105	<b>0.408</b>	0.630	1.718	0.556	5.248	22.84

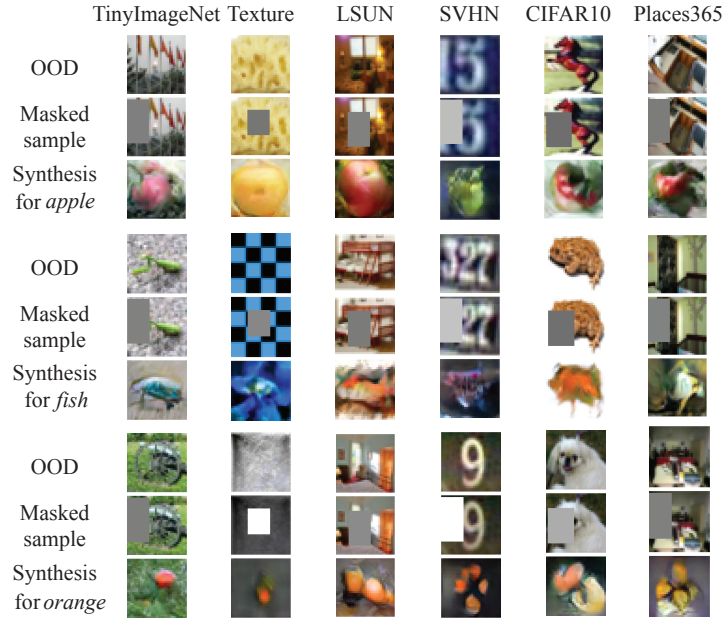


Fig. 4: OOD visualization results of MOODCAT trained on CIFAR-100.

## D.2 Failure Cases

Fig. 5 demonstrates some of MOODCAT’s failure cases. In Fig. 5 (a), the OOD samples sourcing from CIFAR-100, are falsely distinguished as In-D samples (CIFAR-10). As can be seen, OODs and their synthetic images resemble to each other for same degree. For example, the first column’s “cattle” partly contains some features such as legs and the tail, which match the given semantic label “horse” well, resulting in the synthesis having high image quality while resembling to the input image, therefore leading to the final misjudgement.

Fig. 5 (b) presents several False Negative samples, i.e., samples sourcing from In-D are wrongly predicted as OOD samples. As can be observed, the In-D sample with rare characteristics, e.g. a blue fog, an ostrich with its head down, are more likely to be misclassified as OOD. In addition, if the mask happens to cover the object completely, MOODCAT can hardly recover the input image without necessary features, as the cases shown in the third and fourth columns of Fig. 5 (b). Moreover, an poor semantic meaning in the In-D sample itself can lead to the final misclassification. For example, in the last column of Fig. 5 (b), even humans can hardly tell what is depicted in the input image, let alone MOODCAT.

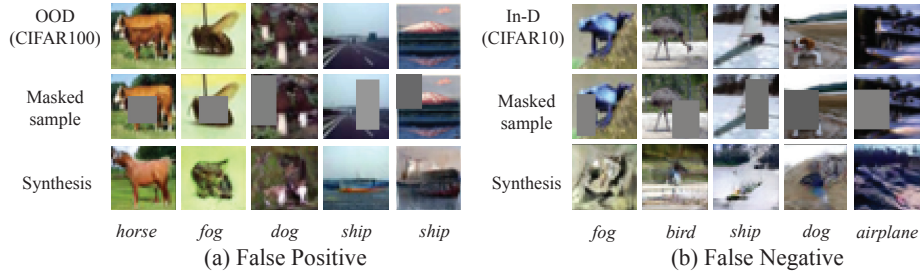


Fig. 5: Failure cases of MOODCAT. We exemplify both False Positive and False Negative failure cases in (a) and (b), respectively. (a) False Positive failure cases, where samples come from OOD dataset (CIFAR-100) are falsely identified as In-D samples (CIFAR-10). (b) False Negative failure cases, where samples belong to In-D are wrongly flagged as OOD samples. The predicted label for each input sample are provided under the corresponding synthetic image.

## References

1. Brock, A., Donahue, J., Simonyan, K.: Large scale gan training for high fidelity natural image synthesis. In: International Conference on Learning Representations (2018)
2. He, K., Chen, X., Xie, S., Li, Y., Dollár, P., Girshick, R.: Masked autoencoders are scalable vision learners. arXiv preprint arXiv:2111.06377 (2021)
3. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: ICLR (Poster) (2015)
4. Schonfeld, E., Schiele, B., Khoreva, A.: A u-net based discriminator for generative adversarial networks. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8207–8216 (2020)