

---

# Category-Extensible Out-of-Distribution Detection via Hierarchical Context Descriptions

---

## Supplementary Materials

### 1 A Implementation Details

#### 2 A.1 Perturbation Guidance

3 In the manuscript, we present a perturbation-guided approach to synthesize spurious samples to train  
4 the hierarchical contexts. Recall that for arbitrary  $k$ -th ID category, we use the spurious context to  
5 explicitly describe a corresponding spurious category, and a critical consideration is how to synthesize  
6 training samples spurious to that  $k$ -th ID category. Recently, generating adversarial data samples  
7 have been widely studied, including GAN networks [19, 14], diffusion models [22, 23], image  
8 attacks [17, 38], and feature-space sampling [7, 26]. For simplicity, we just take the tractable feature-  
9 space sampling as NPOS [26] to generate spurious candidates. In practice, we calculate the  $k$ -NN  
10 distance for each ID sample in the specific category, and generate spurious candidates by sampling  
11 from a multivariate Gaussian distribution around those samples with largest distances (basically away  
12 from the clustering center). Then, we leverage the perturbed descriptions of perceptual context to  
13 guide the further filtering for high-quality spurious syntheses.

14 Given the perceptual context  $\mathbf{v}_k^p = [v_{k,1}^p; v_{k,2}^p; \dots; v_{k,m}^p]$  of  $k$ -th ID category, we randomly apply a  
15 perturbation  $u$  onto one arbitrary  $v_{k,j}^p$  to produce a perturbed description  $\hat{\mathbf{w}}_k^p$  through text-encoder.  
16 Intuitively, there are two ways to perturb a context  $\mathbf{v}_k^p$ : erasing or replacing the specific visual  
17 character  $v_{k,j}^p$ . Consequently, we design three types of perturbation: (1) masking with a placeholder  
18  $u = [\text{MASK}]$ , (2) noise from a Gaussian distribution  $u = \sigma$ , and (3) swapping with another category  
19  $u = v_{k',j'}^p$ . And the perturbed text-feature is produced by:  $\hat{\mathbf{w}}_k^p = \mathcal{T}([v_{k,1}^p; \dots; u; \dots; v_{k,m}^p; \text{CLS}_k])$ .  
20 We also conduct empirical experiments to verify the effectiveness of those perturbations.

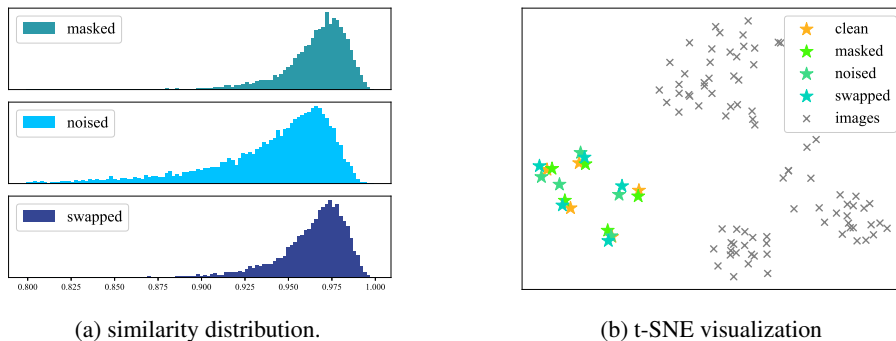


Figure A1: **Statistics of perturbations**, including (a) similarities between original and perturbed text-features, and (b) distribution of original text-feature, perturbed text-features, and image-features.

21 As shown in Fig. A1, all of the perturbed text-features  $\hat{\mathbf{w}}_k^p$  slightly deviate from the original  $\mathbf{w}_k^p$  while  
22 keep the affinity (e.g., shares a 97% similarity against the original one.) Specifically, the noised  $\hat{\mathbf{w}}_k^p$   
23 leads to a greater deviation, since the noised visual character  $v_{k,j}^p := u$  is more unpredictable than  
24 the masked or swapped ones.

25 In addition, now that every perturbation can directly produce the description (*i.e.*, text-feature) of  
 26 an unknown spurious category, one may try to take the perturbed description as a substitute for the  
 27 learned spurious contexts to execute OOD detection. That is, use the perturbed  $\hat{w}_k^p$  to replace the  
 28 learned  $w_k^s$  in Eq.(5) in the manuscript. And the results are shown in Tab. A1, where ImageNet-  
 29 100 [18] is the ID dataset. Given a baseline model [37] learned with perceptual contexts only, simply  
 30 using the perturbed descriptions (denoted as +*Perturb-Desc.*) brings slight improvements (*e.g.*, 0.2%  
 31 decrease on FPR95). The insignificant advantage is not due to the limited capacity of only one  
 32 perturbed description for each ID category. Because ensembling [18] several perturbed descriptions  
 33 for an ID category at once (denoted as +*Perturb-Ensem.*) dose not bring remarkable improvements.  
 34 In contrast, our proposed CATEX can significantly enhance the OOD detection performance, which  
 35 demonstrates it is still necessary to explicitly learn the spurious contexts for each ID category.

Table A1: Comparison with directly using perturbed descriptions for OOD detection.

Method	FPR95↓	AUROC↑
baseline [37]	13.07	97.42
+Perturb-Desc.	12.84	97.43
+Perturb-Ensem.	12.87	97.45
<b>CATEX (Ours)</b>	<b>10.31</b>	<b>97.82</b>

## 36 A.2 Cross-ID-Domain Generalization

37 As indicated in the manuscript, the precise category boundary learned by our method shows robust  
 38 OOD performance when the ID data is shifted. In fact, the shifted ID classification can be further  
 39 boosted by our proposed integrated inference strategy (Eq.(5) in the manuscript), as shown in Tab. A2.  
 40 It implies the regularization item  $\gamma$  successfully modulates the relative similarities between input  
 41 images and learned perceptual descriptions for each category, leading to more precise boundaries.

Table A2: Additionally improved ID accuracy on shifted datasets.

Method	Target Datasets		
	ImageNet-A	ImageNet-R	ImageNet-Sketch
CATEX	50.87	76.67	48.59
+IntegInfer.	<b>50.98</b>	<b>76.72</b>	<b>48.65</b>

42 However, our method only takes the secondary place on ImageNet-Sketch [29] on both ID classifi-  
 43 cation (inferior to NPOS [26]) and OOD detection (inferior to MCM [18]). It is mainly because of  
 44 the huge domain gap between vanilla ImageNet-1K [4] and shifted ImageNet-Sketch. As shown in  
 45 Fig. A2, compared to the shifted ImageNet-A [10] and ImageNet-R [13], images from ImageNet-  
 46 Sketch only preserve objects' shape and main texture, while the color information is totally vanished.  
 47 We leave the generalization to heavily-shifted ID datasets as future work.

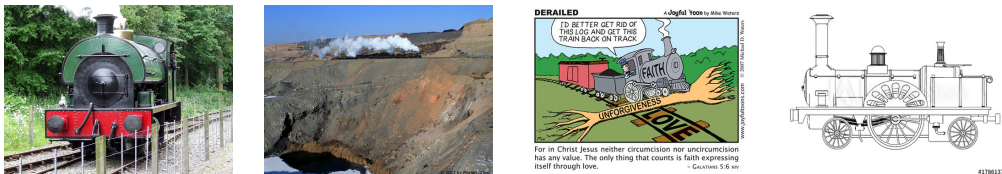


Figure A2: Left to right: examples from ImageNet, ImageNet-A, ImageNet-R, and ImageNet-Sketch.

## 48 A.3 Cross-ID-Task Generalization

49 To verify the efficacy of our proposed framework, we conduct a category-extended experiment in  
 50 Sec.4.1 and Tab.3. Here more implementation details are provided for reproducibility.

51 Given two models independently trained on the separated ImageNet-100 (I) and ImageNet-100 (II),  
 52 how to test them on the union ImageNet-200 ( $I \cup II$ ) with our CATEX is simple. In the vision-  
 53 language prompt-tuning framework, the image-encoder  $\mathcal{I}$  and text-encoder  $\mathcal{T}$  are frozen, and we only  
 54 learn the perceptual and spurious contexts (*i.e.*,  $\mathbf{v}^p$  and  $\mathbf{v}^s$ ). And the  $l_2$ -normalized text-feature can  
 55 be pre-extracted with the 100 category names in each subset, taking the perceptual descriptions for  
 56 example, which are denoted as  $\{\mathbf{w}_{I,k}^p = \mathcal{T}(\mathbf{v}_{I,k}^p; \text{CLS}_{I,k})\}_{k=1}^{100}$  and  $\{\mathbf{w}_{II,k}^p = \mathcal{T}(\mathbf{v}_{II,k}^p; \text{CLS}_{II,k})\}_{k=1}^{100}$ .  
 57 During inference, one may concatenate the 200 text-features together as  $\{\mathbf{w}_k^p\}_{k=1}^{200}$ . Given an input  
 58 image  $I$ , the  $l_2$ -normalized image-feature is extracted by  $\mathbf{x} = \mathcal{I}(I)$ , and the perceptual image-text  
 59 similarities are computed as  $\mathbf{s}^p = [\langle \mathbf{w}_1^p, \mathbf{x} \rangle, \langle \mathbf{w}_2^p, \mathbf{x} \rangle, \dots, \langle \mathbf{w}_{200}^p, \mathbf{x} \rangle] \triangleq [s_1^p, s_2^p, \dots, s_{200}^p]$ . Similarly,  
 60 the spurious similarities become  $\mathbf{s}^s = [s_1^s, s_2^s, \dots, s_{200}^s]$ . Then we can leverage the measurement  
 61 defined in Eq.(5) for both ID classification and OOD detection.

62 As for the competitors, (*e.g.*, VOS [7] and NPOS [26]), two image-encoders are trained separately  
 63 (denoted as  $\mathcal{I}_I$  and  $\mathcal{I}_{II}$ ). And for each input image  $I$ , there are two corresponding image-features:  
 64  $\mathbf{x}_I = \mathcal{I}_I(I)$  and  $\mathbf{x}_{II} = \mathcal{I}_{II}(I)$ . Consequently, there also two sets of image-text similarity vector:  $\mathbf{s}_I =$   
 65  $[\langle \mathbf{w}_1, \mathbf{x}_I \rangle, \langle \mathbf{w}_2, \mathbf{x}_I \rangle, \dots, \langle \mathbf{w}_{200}, \mathbf{x}_I \rangle] = \{\langle \mathbf{w}_k, \mathbf{x}_I \rangle\}_{k=1}^{200}$  and  $\mathbf{s}_{II} = \{\langle \mathbf{w}_k, \mathbf{x}_{II} \rangle\}_{k=1}^{200}$  (the superscript  
 66  $p$  is hidden for simplicity). For compatibility, we choose the one for ID classification and OOD  
 67 detection according to its highest image-text similarity.  $\mathbf{s} = \begin{cases} \mathbf{s}_I & \max(\mathbf{s}_I) > \max(\mathbf{s}_{II}) \\ \mathbf{s}_{II} & \text{otherwise} \end{cases}$ . Now, the  
 68 performance of our method and other rivals are evaluated under the same measurements.

69 Note that since we only take one image encoder throughout, the inference time is fixed (because the  
 70 text-features can be pre-extracted). In contrast, applying other methods brings multiple time cost  
 71 (*e.g.*, twice slower than ours in this case). When the training subsets extend intensely (*e.g.*, from  
 72 ImageNet-1K to ImageNet-21K in our manuscript), our method still keeps a fast speed (*e.g.*, 100FPS  
 73 on V100) during inference, which can even enable real-time applications in practice.

#### 74 A.4 Error Bars

75 To verify the robustness, we repeat the training of our method and the rivals on ImageNet-100 [18]  
 76 with CLIP-B/16 for 3 times, and the results are shown in Tab. A3. Our CATEX consistently  
 77 outperforms the rivals on OOD detection by a significant margin.

Table A3: Error Bars on ImageNet-100 after 3 runs

Method	ACC $\uparrow$	FPR95 $\downarrow$	AUROC $\uparrow$
MSP [8]	<b>94.77</b> ( $\pm 0.05$ )	41.90 ( $\pm 0.61$ )	93.38 ( $\pm 0.05$ )
Energy [16]	<b>94.77</b> ( $\pm 0.05$ )	31.89 ( $\pm 0.50$ )	94.53 ( $\pm 0.18$ )
VOS [7]	94.75 ( $\pm 0.07$ )	24.48 ( $\pm 0.71$ )	96.04 ( $\pm 0.36$ )
NPOS [26]	94.34 ( $\pm 0.12$ )	17.32 ( $\pm 0.87$ )	96.46 ( $\pm 0.13$ )
<b>CATEX</b>	94.11 ( $\pm 0.03$ )	<b>10.97</b> ( $\pm 0.79$ )	<b>97.75</b> ( $\pm 0.07$ )

#### 78 A.5 Software and Hardware

79 We use Python 3.7.13 and PyTorch 1.8.1, and 2 NVIDIA V100-32G GPUs.

## 80 B Experiments on CIFAR Benchmarks

81 To further verify the robustness of our method, we conduct additional experiments on CIFAR-10  
 82 and CIFAR-100 datasets, and evaluate the OOD detection performance on SCOOD [34] benchmark.  
 83 We train our CATEX for 20 epochs, and the other settings are the same as Sec.4 in the manuscript.  
 84 The results are shown in Tab. A4 and Tab. A5, where ‘‘Surr.’’ means the extra TinyImages80M [27]  
 85 is adopted for surrogate OOD training set. Accordingly, our CATEX consistently outperforms the  
 86 competitors as well, and even surpasses those who adopts the extra OOD training data. It implies the  
 87 pre-trained knowledge for large-scale CLIP [20] model leveraged by our method is capable enough  
 88 to detect the OOD samples in the open-world. The efficacy of our CATEX is further demonstrated.

Table A4: Performance on CIFAR-10.

Method	Surr.	ID-ACC $\uparrow$	FPR95 $\downarrow$	AUROC $\uparrow$
MCM [18]	$\times$	90.79	23.14	94.68
ODIN [15]	$\times$	95.36	52.00	82.00
Energy [16]	$\times$	95.36	50.03	83.83
OE [12]	$\checkmark$	94.90	50.53	88.93
UDG [34]	$\checkmark$	94.71	36.22	93.78
CATEX	$\times$	<b>95.57</b>	<b>21.17</b>	<b>95.33</b>

Table A5: Performance on CIFAR-100.

Method	Surr.	ID-ACC $\uparrow$	FPR95 $\downarrow$	AUROC $\uparrow$
MCM [18]	$\times$	66.91	71.93	79.39
ODIN [15]	$\times$	81.84	81.89	77.98
Energy [16]	$\times$	81.84	83.66	79.31
OE [12]	$\checkmark$	81.31	80.06	78.46
UDG [34]	$\checkmark$	80.89	75.45	79.63
CATEX	$\times$	<b>81.99</b>	<b>67.95</b>	<b>84.04</b>

## 90 C Combination with Post-hoc Enhancements

91 Recently, post-hoc OOD detection methods that enhance the single-vision-modal networks (*e.g.*,  
 92 ResNet [9] and ViT [6]) have been widely studied [3, 24, 25, 39, 5]. In this section, we make a step  
 93 towards combining vision-language models with previous post-hoc enhancements for better OOD  
 94 performance. The results are shown in Tab. A6, where ReAct [24] achieves a remarkable improvement.  
 95 It indicates that pruning the extreme feature values according to the unified distributional statistics  
 96 may be more suitable for VLMs to reduce the overconfidence on OOD samples. We hope this can  
 97 bring new insights to the community.

Table A6: Combination with post-hoc methods.

Ccombine	ImageNet-100		ImageNet-1K	
	FPR95 $\downarrow$	AUROC $\uparrow$	FPR95 $\downarrow$	AUROC $\uparrow$
None	10.31	97.82	29.66	93.48
ReAct [24]	<b>10.06</b>	97.82	<b>27.56</b>	<b>93.77</b>
BATS [39]	10.16	<b>97.84</b>	29.37	93.59
ASH [5]	10.19	97.81	29.14	93.27

## 98 D Additional Analysis

99 **Performance improvement.** To further evaluate the improvement brought by our method (*e.g.*,  
 100 8% decrease of FPR95 against NPOS), we conduct a comparative experiment on ImageNet-1K. To  
 101 provide a unified analysis across two models, we take a third-party ResNet-50 model [32] (pre-trained  
 102 on ImageNet-1K classification only) to produce the Maximum SoftMax Probability for each OOD  
 103 sample that is correctly detected by our CATEX while wrongly viewed as ID samples by NPOS.  
 104 According to Fig. A3, our method consistently improves the OOD detection on each interval, where  
 105 the high-probability OOD (generally hard samples) detection is significantly enhanced. It indicates  
 106 that properly leveraging the prior knowledge from pre-trained VLMs can alleviate the OOD problem  
 107 when the fine-tuned visual features are indistinguishable, which is consistent with our motivation.

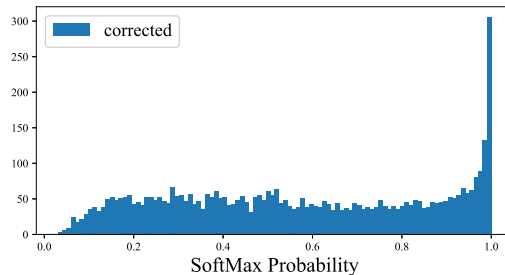


Figure A3: Corrected OOD detections compared with NPOS. The softmax probability predictions on those OOD samples are produced by another pre-trained ResNet-50 [32] classifier.

108 **Failure cases.** As our method still gets 29% FPR95 on ImageNet-1K, we provide some failure cases  
 109 in Fig. A4, which can be summarized into three kinds:

- 110 • Noisy label, where the ID objects (*e.g., dam*) also exists in some OOD images from the test  
111 set. And the dataset composition may need a further examination.
  - 112 • Similar texture, shared by some OOD samples (*e.g., flower*) against ID images (*e.g., starfish*),  
113 and the pre-trained encoders of CLIP are unable to distinguish their features. Applying  
114 image-level spurious OOD syntheses (*e.g., image attacks* [17, 38]) may reduce the texture-  
115 bias.
  - 116 • Same background (*e.g., sky*) that seizes a large proportion of the image may lead to similar  
117 feature representations. Adopting image-level automatic masking techniques [1, 35] to  
118 synthesize spurious OOD samples may alleviate such problem.
- 119 Similar failure cases are also observed in recent SOTA methods, which reveal the unsolved challenges  
120 of OOD detection and suggest the potential directions for future works.

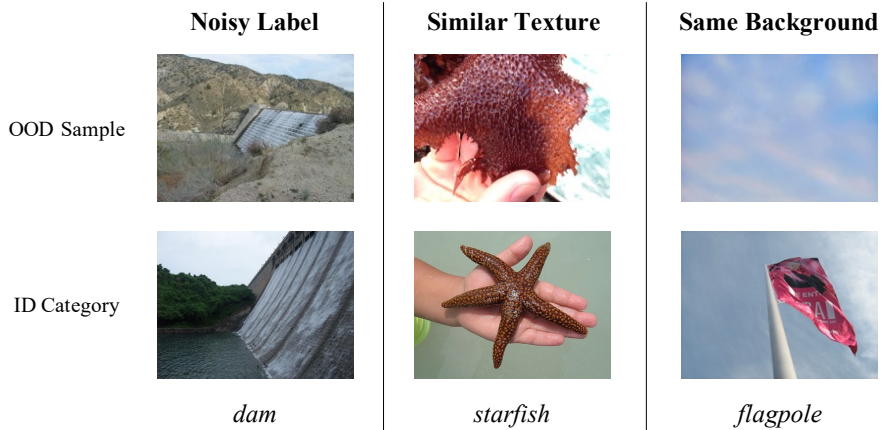


Figure A4: Failed OOD detections of our CATEX

## 121 E Datasets and Baselines

122 For reproducibility, we present the details of datasets and baselines as follows.

123 **ImageNet-100 (I).** Following MCM [18], we take the randomly-sampled 100 classes from ImageNet-  
124 1K [4] as the ImageNet-100 (I) subset, which contains the following categories: n03877845, n03000684,  
125 n03110669, n03710721, n02825657, n02113186, n01817953, n04239074, n02002556, n04356056, n03187595, n03355925, n03125729,  
126 n02058221, n01580077, n03016953, n02843684, n04371430, n01944390, n03887697, n04037443, n02493793, n01518878, n03840681,  
127 n04179913, n01871265, n03866082, n03180011, n01910747, n03388549, n03908714, n01855032, n02134084, n03400231, n04483307,  
128 n03721384, n02033041, n01775062, n02808304, n13052670, n01601694, n04136333, n03272562, n03895866, n03995372, n06785654,  
129 n02111889, n03447721, n03666591, n04376876, n03929855, n02128757, n02326432, n07614500, n01695060, n02484975, n02105412,  
130 n04090263, n03127925, n04550184, n04606251, n02488702, n03404251, n03633091, n02091635, n03457902, n02233338, n02483362,  
131 n04461696, n02871525, n01689811, n01498041, n02107312, n01632458, n03394916, n04147183, n04418357, n03218198, n01917289,  
132 n02102318, n02088364, n09835506, n02095570, n03982430, n04041544, n04562935, n03933933, n01843065, n02128925, n02480495,  
133 n03425413, n03935335, n02971356, n02124075, n07714571, n03133878, n02097130, n02113799, n09399592, n03594945.

134 **ImageNet-100 (II).** Disjoint from ImageNet-100 (I), ImageNet-100 (II) contains another 100 classes  
135 randomly sampled from ImageNet-1K: n02096177, n03769881, n01629819, n04033995, n04357314, n02101388, n02328150,  
136 n03729826, n02655020, n01985128, n02109525, n07715103, n02099429, n04517823, n02088632, n03207743, n03657121, n02948072,  
137 n02106662, n01631663, n09229709, n03793489, n03776460, n07860988, n02129604, n03483316, n02107574, n07716358, n04208210,  
138 n02107908, n04372370, n02119022, n12144580, n01693334, n04548280, n03785016, n03535780, n03599486, n02859443, n04335435,  
139 n02110341, n03902125, n04146614, n01774750, n03314780, n03045698, n01697457, n02869837, n02276258, n04081281, n03956157,  
140 n02487347, n04311174, n02094114, n04409515, n03028079, n03384352, n04532106, n02087394, n04612504, n02100583, n11939491,  
141 n02107142, n01669191, n12998815, n04522168, n02894605, n03529860, n10148035, n01677366, n03775071, n03208938, n04238763,  
142 n02363005, n02804414, n02106382, n03950228, n02128385, n02028035, n04099969, n02481823, n01729322, n02939185, n02483708,  
143 n04162706, n03857828, n02093647, n02927161, n03160309, n02840245, n03920288, n07871810, n04404412, n03947888, n04509417,  
144 n02086910, n02256656, n02412080, n02410509, n03584829.

145 **ImageNet-21K.** The ImageNet-21K dataset on which we conduct the category-extended experiment  
146 is the official winter 2021 released version <sup>1</sup>. For pre-processing, we follow Ridnik *et al* [21] to clean  
147 invalid classes, allocating 50 images per class for validation, and crop-resizing all the images to 224  
148 resolution. Training settings are the same as Sec.4 in our manuscript.

149 **OOD datasets.** Following the literature [30, 26, 31, 18], we mainly consider subsets of  
150 iNaturalist [28], SUN [33], Places [36], and Texture [2] as the OOD datasets, which contains  
151 35640 images in total.

152 **Baselines.** To evaluate the baselines on our experiment settings, we re-implement the most represen-  
153 tative and relevant methods, including MSP [11, 8], Energy [16], VOS [7], and NPOS [26]. For a  
154 fair comparison, we train all the baselines with NPOS’s codebase <sup>2</sup>, and only fine-tune the last two  
155 transformer blocks of image encoder [26].

- 156 • For MSP and Energy, we train a single model with standard cross-entropy loss function for  
157 ID classification only, and infer with respective OOD metrics.
- 158 • For VOS, we take the likelihood-based sampling strategy to generate spurious OOD synthe-  
159 ses, and train the model with uncertainty regularization as suggested [7].
- 160 • For NPOS, we take the non-parametric distance-based sampling strategy to generate spurious  
161 OOD syntheses, and train the model with open-set ERM as suggested [26].

## 162 References

- 163 [1] Daniel Bolya, Cheng-Yang Fu, Xiaoliang Dai, Peizhao Zhang, Christoph Feichtenhofer, and Judy Hoffman.  
164 Token merging: Your ViT but faster. In *International Conference on Learning Representations*, 2023.
- 165 [2] Mircea Cimpoi, Subhransu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing  
166 textures in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
167 pages 3606–3613, 2014.
- 168 [3] Jeremy Cohen, Elan Rosenfeld, and Zico Kolter. Certified adversarial robustness via randomized smoothing.  
169 In *international conference on machine learning*, pages 1310–1320. PMLR, 2019.
- 170 [4] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical  
171 image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255.  
172 IEEE, 2009.
- 173 [5] Andrija Djuricic, Nebojsa Bozanic, Arjun Ashok, and Rosanne Liu. Extremely simple activation shaping  
174 for out-of-distribution detection. In *International Conference on Learning Representations*, 2023.
- 175 [6] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas  
176 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and  
177 Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. *International  
178 Conference on Machine Learning*, 2021.
- 179 [7] Xuefeng Du, Zhaoning Wang, Mu Cai, and Yixuan Li. Vos: Learning what you don’t know by virtual  
180 outlier synthesis. In *International Conference on Learning Representations*, 2022.
- 181 [8] Stanislav Fort, Jie Ren, and Balaji Lakshminarayanan. Exploring the limits of out-of-distribution detection.  
182 *Advances in Neural Information Processing Systems*, 34:7068–7081, 2021.
- 183 [9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.  
184 In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- 185 [10] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai,  
186 Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical analysis of  
187 out-of-distribution generalization. In *Proceedings of the IEEE/CVF International Conference on Computer  
188 Vision*, pages 8340–8349, 2021.
- 189 [11] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples  
190 in neural networks. In *International Conference on Learning Representations*, 2017.

---

<sup>1</sup><https://image-net.org/>

<sup>2</sup><https://github.com/deeplearning-wisc/npos>

- 191 [12] Dan Hendrycks, Mantas Mazeika, and Thomas Dietterich. Deep anomaly detection with outlier exposure.  
192 *Proceedings of the International Conference on Learning Representations*, 2019.
- 193 [13] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial  
194 examples. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,  
195 pages 15262–15271, 2021.
- 196 [14] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial  
197 networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages  
198 4401–4410, 2019.
- 199 [15] Shiyu Liang, Yixuan Li, and R Srikant. Enhancing the reliability of out-of-distribution image detection in  
200 neural networks. In *International Conference on Learning Representations*, 2018.
- 201 [16] Weitang Liu, Xiaoyun Wang, John Owens, and Yixuan Li. Energy-based out-of-distribution detection.  
202 *Advances in neural information processing systems*, 33:21464–21475, 2020.
- 203 [17] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. To-  
204 wards deep learning models resistant to adversarial attacks. In *International Conference on Learning*  
205 *Representations*, 2018.
- 206 [18] Yifei Ming, Ziyang Cai, Jiuxiang Gu, Yiyu Sun, Wei Li, and Yixuan Li. Delving into out-of-distribution  
207 detection with vision-language representations. In *Advances in Neural Information Processing Systems*,  
208 2022.
- 209 [19] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*,  
210 2014.
- 211 [20] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish  
212 Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from  
213 natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR,  
214 2021.
- 215 [21] Tal Ridnik, Emanuel Ben-Baruch, Asaf Noy, and Lihl Zelnik-Manor. Imagenet-21k pretraining for the  
216 masses. *arXiv preprint arXiv:2104.10972*, 2021.
- 217 [22] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution  
218 image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer*  
219 *Vision and Pattern Recognition*, pages 10684–10695, 2022.
- 220 [23] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar  
221 Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-  
222 image diffusion models with deep language understanding. *Advances in Neural Information Processing*  
223 *Systems*, 35:36479–36494, 2022.
- 224 [24] Yiyu Sun, Chuan Guo, and Yixuan Li. React: Out-of-distribution detection with rectified activations. In  
225 *Advances in Neural Information Processing Systems*, 2021.
- 226 [25] Yiyu Sun and Yixuan Li. Dice: Leveraging sparsification for out-of-distribution detection. In *Computer*  
227 *Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part*  
228 *XXIV*, pages 691–708. Springer, 2022.
- 229 [26] Leitian Tao, Xuefeng Du, Jerry Zhu, and Yixuan Li. Non-parametric outlier synthesis. In *The Eleventh*  
230 *International Conference on Learning Representations*, 2023.
- 231 [27] Antonio Torralba, Rob Fergus, and William T Freeman. 80 million tiny images: A large data set  
232 for nonparametric object and scene recognition. *IEEE transactions on pattern analysis and machine*  
233 *intelligence*, 30(11):1958–1970, 2008.
- 234 [28] Grant Van Horn, Oisin Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam, Pietro  
235 Perona, and Serge Belongie. The inaturalist species classification and detection dataset. In *Proceedings of*  
236 *the IEEE conference on computer vision and pattern recognition*, pages 8769–8778, 2018.
- 237 [29] Haohan Wang, Songwei Ge, Zachary Lipton, and Eric P Xing. Learning robust global representations  
238 by penalizing local predictive power. In *Advances in Neural Information Processing Systems*, pages  
239 10506–10518, 2019.
- 240 [30] Haotao Wang, Aston Zhang, Yi Zhu, Shuai Zheng, Mu Li, Alex J Smola, and Zhangyang Wang. Partial and  
241 asymmetric contrastive learning for out-of-distribution detection in long-tailed recognition. In *International*  
242 *Conference on Machine Learning*, pages 23446–23458. PMLR, 2022.

- 243 [31] Qizhou Wang, Junjie Ye, Feng Liu, Quanyu Dai, Marcus Kalander, Tongliang Liu, HAO Jianye, and  
244 Bo Han. Out-of-distribution detection with implicit outlier transformation. In *The Eleventh International  
245 Conference on Learning Representations*, 2023.
- 246 [32] Ross Wightman. Pytorch image models. <https://github.com/rwightman/pytorch-image-models>,  
247 2019.
- 248 [33] Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database: Large-  
249 scale scene recognition from abbey to zoo. In *2010 IEEE computer society conference on computer vision  
250 and pattern recognition*, pages 3485–3492. IEEE, 2010.
- 251 [34] Jingkan Yang, Haoqi Wang, Litong Feng, Xiaopeng Yan, Huabin Zheng, Wayne Zhang, and Ziwei  
252 Liu. Semantically coherent out-of-distribution detection. In *Proceedings of the IEEE/CVF International  
253 Conference on Computer Vision*, pages 8301–8309, 2021.
- 254 [35] Tao Yu, Runseng Feng, Ruoyu Feng, Jinming Liu, Xin Jin, Wenjun Zeng, and Zhibo Chen. Inpaint  
255 anything: Segment anything meets image inpainting. *arXiv preprint arXiv:2304.06790*, 2023.
- 256 [36] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million  
257 image database for scene recognition. *IEEE transactions on pattern analysis and machine intelligence*,  
258 40(6):1452–1464, 2017.
- 259 [37] Kaiyang Zhou, Jingkan Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language  
260 models. *International Journal of Computer Vision*, 130(9):2337–2348, 2022.
- 261 [38] Yao Zhu, Yuefeng Chen, Xiaodan Li, Kejiang Chen, Yuan He, Xiang Tian, Bolun Zheng, Yaowu Chen,  
262 and Qingming Huang. Toward understanding and boosting adversarial transferability from a distribution  
263 perspective. *IEEE Transactions on Image Processing*, 31:6487–6501, 2022.
- 264 [39] Yao Zhu, YueFeng Chen, Chuanlong Xie, Xiaodan Li, Rong Zhang, Xiang Tian, Yaowu Chen, et al.  
265 Boosting out-of-distribution detection with typical features. In *Advances in Neural Information Processing  
266 Systems*, 2022.