



DICE: Leveraging Sparsification for Out-of-Distribution Detection

Yiyou Sun  and Yixuan Li 

Computer Science Department, University of Wisconsin-Madison
{sunyiyou, sharonli}@cs.wisc.edu

Abstract. Detecting out-of-distribution (OOD) inputs is a central challenge for safely deploying machine learning models in the real world. Previous methods commonly rely on an OOD score derived from the overparameterized weight space, while largely overlooking the role of *sparsification*. In this paper, we reveal important insights that reliance on unimportant weights and units can directly attribute to the brittleness of OOD detection. To mitigate the issue, we propose a sparsification-based OOD detection framework termed **DICE**. Our key idea is to rank weights based on a measure of contribution, and selectively use the most salient weights to derive the output for OOD detection. We provide both empirical and theoretical insights, characterizing and explaining the mechanism by which DICE improves OOD detection. By pruning away noisy signals, DICE provably reduces the output variance for OOD data, resulting in a sharper output distribution and stronger separability from ID data. We demonstrate the effectiveness of sparsification-based OOD detection on several benchmarks and establish competitive performance. Code is available at: <https://github.com/deeplearning-wisc/dice.git>.

Keywords: Out-of-distribution Detection, Sparsification

1 Introduction

Deep neural networks deployed in real-world systems often encounter out-of-distribution (OOD) inputs—samples from unknown classes that the network has not been exposed to during training, and therefore should not be predicted by the model in testing. Being able to estimate and mitigate OOD uncertainty is paramount for safety-critical applications such as medical diagnosis [48, 60] and autonomous driving [10]. For example, an autonomous vehicle may fail to recognize objects on the road that do not appear in its detection model’s training set, potentially leading to a crash. This gives rise to the importance of OOD detection, which allows the learner to express ignorance and take precautions in the presence of OOD data.

The main challenge in OOD detection stems from the fact that modern deep neural networks can easily produce overconfident predictions on OOD inputs, making the separation between in-distribution (ID) and OOD data a non-trivial task. The vulnerability of machine learning to OOD data can be hard-wired in

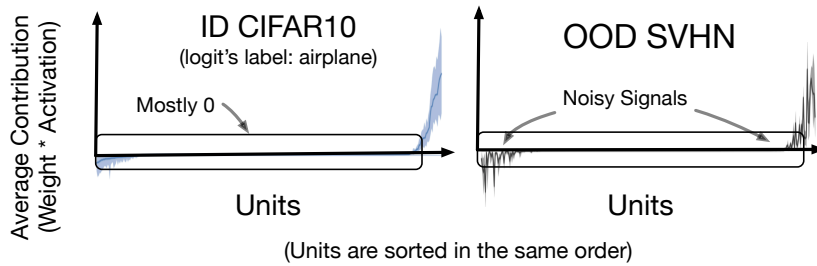


Fig. 1. Illustration of unit contribution (*i.e.*, $\text{weight} \times \text{activation}$) to the class output. For class c , the output $f_c(\mathbf{x})$ is the summation of unit contribution from the penultimate feature layer of a neural network. *Units are sorted in the same order*, based on the expectation of ID data’s contribution (averaged over many CIFAR-10 samples) on the x -axis. **Shades indicate the variance for each unit.** **Left:** For in-distribution data (CIFAR-10, airplane), only a subset of units contributes to the model output. **Right:** In contrast, out-of-distribution (OOD) data can trigger a non-negligible fraction of units with noisy signals, as indicated by the variances.

high-capacity models used in practice. In particular, modern deep neural networks can overfit observed patterns in the training data [67], and worse, activate features on unfamiliar inputs [46]. To date, existing OOD detection methods commonly derive OOD scores using overparameterized weights, while largely overlooking the role of *sparsification*. This paper aims to bridge the gap.

In this paper, we start by revealing key insights that reliance on unimportant units and weights can directly attribute to the brittleness of OOD detection. Empirically on a network trained with CIFAR-10, we show that an OOD image can activate a non-negligible fraction of units in the penultimate layer (see Figure 1, right). Each point on the horizontal axis corresponds to a single unit. The y-axis measures the unit contribution (*i.e.*, $\text{weight} \times \text{activation}$) to the output of class AIRPLANE, with the solid line and the shaded area indicating the mean and variance, respectively. Noticeably, for OOD data (gray), we observe a non-negligible fraction of “noisy” units that display high variances of contribution, which is then aggregated to the model’s output through summation. As a result, such noisy signals can undesirably manifest in model output—increasing the variance of output distribution and reducing the separability from ID data.

The above observation motivates a simple and effective method, *Directed Sparsification (DICE)*, for OOD detection. DICE leverages the observation that a model’s prediction for an ID class depends on only a subset of important units (and corresponding weights), as evidenced in Figure 1 (left). To exploit this, our novel idea is to rank weights based on the measure of contribution, and selectively use the most contributing weights to derive the output for OOD detection. As a result of the weight sparsification, we show that the model’s output becomes more separable between ID and OOD data. Importantly, DICE can be conveniently used by post hoc weight masking on a pre-trained network

and therefore can preserve the ID classification accuracy. Orthogonal to existing works on sparsification for accelerating computation, our primary goal is to explore the sparsification approach for improved OOD detection performance.

We provide both empirical and theoretical insights characterizing and explaining the mechanism by which DICE improves OOD detection. We perform extensive evaluations and establish competitive performance on common OOD detection benchmarks, including CIFAR-10, CIFAR-100 [29], and a large-scale ImageNet benchmark [25]. Compared to the competitive post hoc method ReAct [51], DICE reduces the FPR95 by up to 12.55%. Moreover, we perform ablation using various sparsification techniques and demonstrate the benefit of directed sparsification for OOD detection. Theoretically, by pruning away noisy signals from unimportant units and weights, DICE *provably reduces the output variance* and results in a sharper output distribution (see Section 6). The sharper distributions lead to a stronger separability between ID and OOD data and overall improved OOD detection performance (*c.f.* Figure 2). Our **key results and contributions** are:

- (Methodology) We introduce DICE, a simple and effective approach for OOD detection utilizing post hoc weight sparsification. To the best of our knowledge, DICE is the first to explore and demonstrate the effectiveness of sparsification for OOD detection.
- (Experiments) We extensively evaluate DICE on common benchmarks and establish competitive performance among post hoc OOD detection baselines. DICE outperforms the strong baseline [51] by reducing the FPR95 by up to 12.55%. We show DICE can effectively improve OOD detection while preserving the classification accuracy on ID data.
- (Theory and ablations) We provide ablation and theoretical analysis that improves understanding of a sparsification-based method for OOD detection. Our analysis reveals an important variance reduction effect, which provably explains the effectiveness of DICE. We hope our insights inspire future research on weight sparsification for OOD detection.

2 Preliminaries

We start by recalling the general setting of the supervised learning problem. We denote by $\mathcal{X} = \mathbb{R}^d$ the input space and $\mathcal{Y} = \{1, 2, \dots, C\}$ the output space. A learner is given access to a set of training data $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ drawn from an unknown joint data distribution \mathcal{P} defined on $\mathcal{X} \times \mathcal{Y}$. Furthermore, let \mathcal{P}_{in} denote the marginal probability distribution on \mathcal{X} .

Out-of-distribution detection When deploying a model in the real world, a reliable classifier should not only accurately classify known in-distribution (ID) samples, but also identify any OOD input as “unknown”. This can be achieved through having dual objectives: ID/OOD classification and multi-class classification of ID data [3].

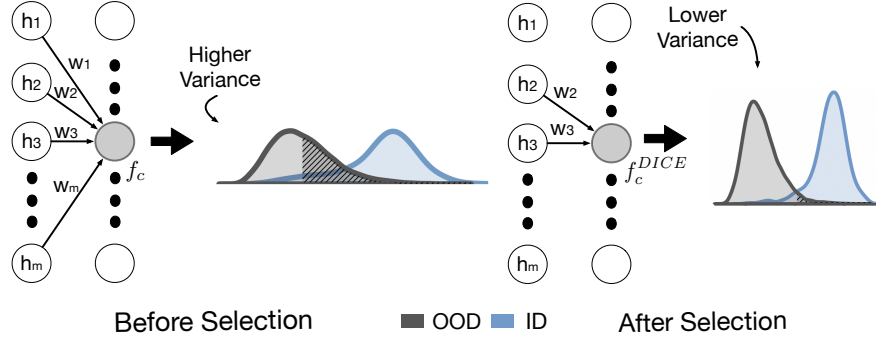


Fig. 2. Illustration of out-of-distribution detection using *Directed Sparsification (DICE)*. We consider a pre-trained neural network, which encodes an input \mathbf{x} to a feature vector $h(\mathbf{x}) \in \mathbb{R}^m$. **Left:** The logit output $f_c(\mathbf{x})$ of class c is a linear combination of activation from *all* units in the preceding layer, weighted by w_i . The full connection results in a high variance for OOD data’s output, as depicted in the gray. **Right:** Our proposed approach leverages a selective subset of weights, which effectively reduces the output variance for OOD data, resulting in a sharper score distribution and stronger separability from ID data. The output distributions are based on CIFAR-10 trained network, with ID class label “frog” and SVHN as OOD.

OOD detection can be formulated as a binary classification problem. At test time, the goal of OOD detection is to decide whether a sample $\mathbf{x} \in \mathcal{X}$ is from \mathcal{P}_{in} (ID) or not (OOD). In literature, OOD distribution \mathcal{P}_{out} often simulates unknowns encountered during deployment time, such as samples from an irrelevant distribution whose label set has no intersection with \mathcal{Y} and therefore should not be predicted by the model. The decision can be made via a thresholding comparison:

$$g_\lambda(\mathbf{x}) = \begin{cases} \text{in} & S(\mathbf{x}) \geq \lambda \\ \text{out} & S(\mathbf{x}) < \lambda \end{cases},$$

where samples with higher scores $S(\mathbf{x})$ are classified as ID and vice versa, and λ is the threshold.

3 Method

Method overview Our novel idea is to selectively use a subset of important weights to derive the output for OOD detection. By utilizing sparsification, the network prevents adding irrelevant information to the output. We illustrate our idea in Figure 2. Without DICE (*left*), the final output is a summation of weighted activations across all units, which can have a high variance for OOD data (colored in gray). In contrast, with DICE (*right*), the variance of output can

be significantly reduced, which improves separability from ID data. We proceed with describing our method in details, and provide the theoretical explanation later in Section 6.

3.1 DICE: Directed Sparsification

We consider a deep neural network parameterized by θ , which encodes an input $\mathbf{x} \in \mathbb{R}^d$ to a feature space with dimension m . We denote by $h(\mathbf{x}) \in \mathbb{R}^m$ the feature vector from the penultimate layer of the network. A weight matrix $\mathbf{W} \in \mathbb{R}^{m \times C}$ connects the feature $h(\mathbf{x})$ to the output $f(\mathbf{x})$.

Contribution matrix We perform a *directed sparsification* based on a measure of contribution, and preserve the most important weights in \mathbf{W} . To measure the contribution, we define a contribution matrix $\mathbf{V} \in \mathbb{R}^{m \times C}$, where each column $\mathbf{v}_c \in \mathbb{R}^m$ is given by:

$$\mathbf{v}_c = \mathbb{E}_{\mathbf{x} \in \mathcal{D}}[\mathbf{w}_c \odot h(\mathbf{x})], \quad (1)$$

where \odot indicates the element-wise multiplication, and \mathbf{w}_c indicates weight vector for class c . Each element in $\mathbf{v}_c \in \mathbb{R}^m$ intuitively measures the corresponding unit's average contribution to class c , estimated empirically on in-distribution data \mathcal{D} . A larger value indicates a higher contribution to the output $f_c(\mathbf{x})$ of class c . The vector \mathbf{v}_c is derived for all classes $c \in \{1, 2, \dots, C\}$, forming the contribution matrix \mathbf{V} . Each element $\mathbf{v}_c^i \in \mathbf{V}$ measures the average contribution (**weight** \times **activation**) from a unit i to the output class $c \in \{1, 2, \dots, C\}$.

We can now select the top- k weights based on the k -largest elements in \mathbf{V} . In particular, we define a masking matrix $\mathbf{M} \in \mathbb{R}^{m \times C}$, which returns a matrix by setting 1 for entries corresponding to the k largest elements in \mathbf{V} and setting other elements to 0. The model output under *contribution-directed sparsification* is given by

$$f^{\text{DICE}}(\mathbf{x}; \theta) = (\mathbf{M} \odot \mathbf{W})^\top h(\mathbf{x}) + \mathbf{b}, \quad (2)$$

where $\mathbf{b} \in \mathbb{R}^C$ is the bias vector. The procedure described above essentially accounts for information from the most relevant units in the penultimate layer. Importantly, the sparsification can be conveniently imposed by *post hoc* weight masking on the final layer of a pre-trained network, without changing any parameterizing of the neural network. Therefore one can improve OOD detection while preserving the ID classification accuracy.

Sparsity parameter p To align with the convention in literature, we use the sparsity parameter $p = 1 - \frac{k}{m \cdot C}$ in the remainder paper. A higher p indicates a larger fraction of weights dropped. When $p = 0$, the output becomes equivalent to the original output $f(\mathbf{x}; \theta)$ using dense transformation, where $f(\mathbf{x}; \theta) = \mathbf{W}^\top h(\mathbf{x}) + \mathbf{b}$. We provide ablations on the sparsity parameter later in Section 5.

3.2 OOD Detection with DICE

Our method DICE in Section 3.1 can be flexibly leveraged by the downstream OOD scoring function:

$$g_\lambda(\mathbf{x}) = \begin{cases} \text{in} & S_\theta(\mathbf{x}) \geq \lambda \\ \text{out} & S_\theta(\mathbf{x}) < \lambda \end{cases}, \quad (3)$$

where a thresholding mechanism is exercised to distinguish between ID and OOD during test time. The threshold λ is typically chosen so that a high fraction of ID data (*e.g.*, 95%) is correctly classified. Following recent work by Liu *et. al* [36], we derive an energy score using the logit output $f^{\text{DICE}}(\mathbf{x}; \theta)$ with contribution-directed sparsification. The function maps the logit outputs $f^{\text{DICE}}(\mathbf{x}; \theta)$ to a scalar $E_\theta(\mathbf{x}) \in \mathbb{R}$, which is relatively lower for ID data:

$$S_\theta(\mathbf{x}) = -E_\theta(\mathbf{x}) = \log \sum_{c=1}^C \exp(f_c^{\text{DICE}}(\mathbf{x}; \theta)). \quad (4)$$

The energy score can be viewed as the log of the denominator in softmax function:

$$p(y|\mathbf{x}) = \frac{p(\mathbf{x}, y)}{p(\mathbf{x})} = \frac{\exp(f_y(\mathbf{x}; \theta))}{\sum_{c=1}^C \exp(f_c(\mathbf{x}; \theta))}, \quad (5)$$

and enjoys better theoretical interpretation than using posterior probability $p(y|\mathbf{x})$. Note that DICE can also be compatible with an alternative scoring function such as maximum softmax probability (MSP) [20], though the performance of MSP is less competitive (see Appendix F). Later in Section 6, we formally characterize and explain why DICE improves the separability of the scores between ID and OOD data.

4 Experiments

In this section, we evaluate our method on a suite of OOD detection tasks. We begin with the CIFAR benchmarks that are routinely used in literature (Section 4.1). In Section 4.2 we continue with a large-scale OOD detection task based on ImageNet.

4.1 Evaluation on Common Benchmarks

Experimental details We use CIFAR-10 [29], and CIFAR-100 [29] datasets as in-distribution data. We use the standard split with 50,000 training images and 10,000 test images. We evaluate the model on six common OOD benchmark datasets: Textures [6], SVHN [45], Places365 [68], LSUN-Crop [66], LSUN-Resize [66], and iSUN [64]. We use DenseNet-101 architecture [23] and train on in-distribution datasets. The feature dimension of the penultimate layer is 342. For both CIFAR-10 and CIFAR-100, the model is trained for 100 epochs

Table 1. Comparison with competitive *post hoc* out-of-distribution detection method on CIFAR benchmarks. All values are percentages and are averaged over 6 OOD test datasets. The full results for each evaluation dataset are provided in Appendix G. We report standard deviations estimated across 5 independent runs. [§] indicates an exception, where model retraining using a different loss function is required.

Method	CIFAR-10		CIFAR-100	
	FPR95 ↓	AUROC ↑	FPR95 ↓	AUROC ↑
MSP [20]	48.73	92.46	80.13	74.36
ODIN [34]	24.57	93.71	58.14	84.49
GODIN [§] [22]	34.25	90.61	52.87	85.24
Mahalanobis [32]	31.42	89.15	55.37	82.73
Energy [36]	26.55	94.57	68.45	81.19
ReAct [51]	26.45	94.95	62.27	84.47
DICE (ours)	20.83± 1.58	95.24± 0.24	49.72± 1.69	87.23± 0.73

with batch size 64, weight decay 0.0001 and momentum 0.9. The start learning rate is 0.1 and decays by a factor of 10 at epochs 50, 75, and 90. We use the validation strategy in Appendix C to select p .

DICE vs. competitive baselines We show the results in Table 1, where DICE outperforms competitive baselines. In particular, we compare with Maximum Softmax Probability [20], ODIN [34], Mahalanobis distance [32], Generalized ODIN [22], Energy score [36], and ReAct [51]. For a fair comparison, all the methods derive the OOD score post hoc from the same pre-trained model, except for G-ODIN which requires model re-training. For readers’ convenience, a brief introduction of baselines and hyperparameters is provided in Appendix B.

On CIFAR-100, we show that DICE reduces the average FPR95 by **18.73%** compared to the vanilla energy score [36] without sparsification. Moreover, our method also outperforms a competitive method ReAct [51] by 12.55%. While ReAct only considers activation space, DICE examines *both the weights and activation* values together—the multiplication of which directly determines the network’s logit output. Overall our method is more generally applicable, and can be implemented through a simple post hoc weight masking.

ID classification accuracy Given the *post hoc* nature of DICE, once the input image is marked as ID, one can always use the original fc layer, which is guaranteed to give identical classification accuracy. This incurs minimal overhead and results in optimal performance for both classification and OOD detection. We also measure the classification accuracy under different sparsification parameter p . Due to the space limit, the full results are available in Table 6 in Appendix.

Table 2. Main results. Comparison with competitive *post hoc* out-of-distribution detection methods. All methods are based on a discriminative model trained on ImageNet. \uparrow indicates larger values are better and \downarrow indicates smaller values are better. All values are percentages. **Bold** numbers are superior results.

Methods	OOD Datasets								Average	
	iNaturalist		SUN		Places		Textures			
	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC
	\downarrow	\uparrow	\downarrow	\uparrow	\downarrow	\uparrow	\downarrow	\uparrow	\downarrow	\uparrow
MSP [20]	54.99	87.74	70.83	80.86	73.99	79.76	68.00	79.61	66.95	81.99
ODIN [34]	47.66	89.66	60.15	84.59	67.89	81.78	50.23	85.62	56.48	85.41
GODIN [22]	61.91	85.40	60.83	85.60	63.70	83.81	77.85	73.27	66.07	82.02
Mahalanobis [32]	97.00	52.65	98.50	42.41	98.40	41.79	55.80	85.01	87.43	55.47
Energy [36]	55.72	89.95	59.26	85.89	64.92	82.86	53.72	85.99	58.41	86.17
ReAct [51]	20.38	96.22	24.20	94.20	33.85	91.58	47.30	89.80	31.43	92.95
DICE (ours)	25.63	94.49	35.15	90.83	46.49	87.48	31.72	90.30	34.75	90.77
DICE + ReAct (ours)	18.64	96.24	25.45	93.94	36.86	90.67	28.07	92.74	27.25	93.40

4.2 Evaluation on ImageNet

Dataset We then evaluate DICE on a large-scale ImageNet classification model. Following MOS [25], we use four OOD test datasets from (subsets of) Places365 [68], Textures [6], iNaturalist [57], and SUN [63] with non-overlapping categories *w.r.t.* ImageNet. The evaluations span a diverse range of domains including fine-grained images, scene images, and textural images. OOD detection for the ImageNet model is more challenging due to both a larger feature space ($m = 2,048$) as well as a larger label space ($C = 1,000$). In particular, the large-scale evaluation can be relevant to real-world applications, where the deployed models often operate on images that have high resolution and contain many class labels. Moreover, as the number of feature dimensions increases, noisy signals may increase accordingly, which can make OOD detection more challenging.

Experimental details We use a pre-trained ResNet-50 model [17] for ImageNet-1k provided by Pytorch. At test time, all images are resized to 224×224 . We use the entire training dataset to estimate the contribution matrix and masking matrix \mathbf{M} . We use the validation strategy in Appendix C to select p . The hardware used for experiments is specified in Appendix A.

Comparison with baselines In Table 2, we compare DICE with competitive post hoc OOD detection methods. We report performance for each OOD test dataset, as well as the average of the four. We first contrast DICE with energy score [36], which allows us to see the direct benefit of using sparsification under the same scoring function. DICE reduces the FPR95 drastically from 58.41% to 34.75%, a **23.66%** improvement using sparsification. Second, we contrast with a recent method ReAct [51], which demonstrates strong performance on this challenging task using activation truncation. With the truncated activation proposed in ReAct [51], we show that DICE can further reduce the FPR95 by 5.78% with weight sparsification. Since the comparison is conducted on the same

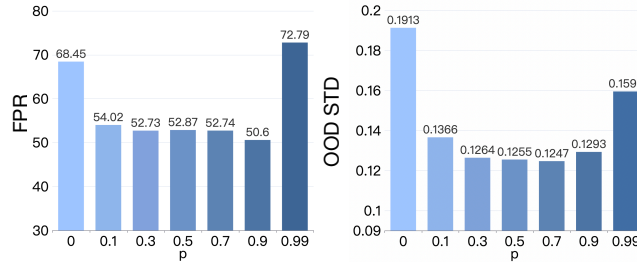


Fig. 3. Effect of varying sparsity parameter p during inference time. Model is trained on CIFAR-100 using DenseNet101 [23].

scoring function and feature activation, the performance improvement from ReAct to DICE+ReAct precisely highlights the benefit of using weight sparsification as opposed to the full weights. Lastly, Mahalanobis displays limiting performance on ImageNet, while being computationally expensive due to estimating the inverse of the covariance matrix. In contrast, DICE is easy to use in practice, and can be implemented through simple post hoc weight masking.

5 Discussion and Ablations

Ablation on sparsity parameter p We now characterize the effect of sparsity parameter p . In Figure 3, we summarize the OOD detection performance for DenseNet trained on CIFAR-100, where we vary $p = \{0.1, 0.3, 0.5, 0.7, 0.9, 0.99\}$. Interestingly, we observe the performance improves with mild sparsity parameter p . A significant improvement can be observed from $p = 0$ (no sparsity) to $p = 0.1$. As we will theoretically later in Section 6, this is because the leftmost part of units being pruned has larger variances for OOD data (gray shade). Units in the middle part have small variances and contributions for both ID and OOD, therefore leading to similar performance as p increases mildly. This ablation confirms that over-parameterization does compromise the OOD detection ability, and DICE can effectively alleviate the problem. In the extreme case when p is too large (*e.g.*, $p = 0.99$), the OOD performance starts to degrade as expected.

Effect of variance reduction for output distribution Figure 2 shows that DICE has an interesting variance reduction effect on the output distribution for OOD data, and at the same time preserves the information for the ID data (CIFAR-10, class “frog”). The output distribution without any sparsity ($p = 0$) appears to have a larger variance, resulting in less separability from ID data (see left of Figure 2). In contrast, sparsification with DICE results in a sharper distribution, which benefits OOD detection. In Figure 3, we also measure the standard deviation of energy score for OOD data (normalized by the mean of ID data’s OOD scores in each setting). By way of sparsification, DICE can reduce the output variance. In Section 6, we formally characterize this and provide a theoretical explanation.

Table 3. Ablation results. Effect of different *post hoc* sparsification methods for OOD detection with ImageNet as ID dataset. All sparsification methods are based on the same OOD scoring function [36], with sparsity parameter $p = 0.7$. All values are percentages and are averaged over multiple OOD test datasets.

Method	FPR95↓	AUROC↑
Weight-Dropout	76.28	76.55
Unit-Dropout	83.91	64.98
Weight-Pruning	52.81	87.08
Unit-Pruning	90.80	49.15
DICE (Ours)	34.75	90.77

Ablation on pruning methods In this ablation, we evaluate OOD detection performance under the most common *post hoc* sparsification methods. Here we primarily consider post hoc sparsification strategy which operates conveniently on a *pre-trained* network, instead of training with sparse regularization or architecture modification. The property is especially desirable for the adoption of OOD detection methods in real-world production environments, where the overhead cost of retraining can be sometimes prohibitive. Orthogonal to existing works on sparsification, our primary goal is to explore the role of sparsification for improved OOD detection performance, rather than establishing a generic sparsification algorithm. We consider the most common strategies, covering both unit-based and weight-based sparsification methods: (1) unit dropout [50] which randomly drops a fraction of units, (2) unit pruning [33] which drops units with the smallest L_2 norm of the corresponding weight vectors, (3) weight dropout [58] which randomly drops weights in the fully connected layer, and (4) weight pruning [16] drops weights with the smallest entries under the L_1 norm. For consistency, we use the same OOD scoring function and the same sparsity parameter for all.

Our ablation reveals several important insights shown in Table 3. First, in contrasting weight dropout vs. DICE, a salient performance gap of 41.53% (FPR95) is observed under the same sparsity. This suggests the importance of dropping weights *directedly* rather than *randomly*. Second, DICE outperforms a popular L_1 -norm-based pruning method [16] by up to 18.06% (FPR95). While it prunes weights with low magnitude, negative weights with large L_1 -norm can be kept. The negative weights can undesirably corrupt the output with noisy signals (as shown in Figure 1). The performance gain of DICE over [16] attributes to our contribution-directed sparsification, which is better suited for OOD detection.

Ablation on unit selection We have shown that choosing a subset of weights (with *top-k* unit contribution) significantly improves the OOD detection performance. In this ablation, we also analyze those “lower contribution units” for OOD detection. Specifically, we consider: (1) *Bottom-k* which only includes k unit contribution with least contribution values, (2) *top+bottom-k* which includes k unit contribution with largest and smallest contribution values, (3) *random-k*

Table 4. Ablation on different strategies of choosing a subset of units. Values are FPR95 (averaged over multiple test datasets).

Method	CIFAR-10↓	CIFAR-100 ↓
Bottom- k	91.87	99.70
(Top+Bottom)- k	24.25	59.93
Random- k	62.12	77.48
Top- k (DICE)	20.83 ± 1.58	49.72 ± 1.69

which randomly includes k unit contribution and (4) *top- k* which is equivalent to DICE method. In Table 4, we show that DICE outperforms these variants.

6 Why does DICE improve OOD detection?

In this section, we formally explain the mechanism by which reliance on irrelevant units hurts OOD detection and how DICE effectively mitigates the issue. Our analysis highlights that DICE reduces the output variance for both ID and OOD data. Below we provide details.

Setup For a class c , we consider the unit contribution vector \mathbf{v} , the element-wise multiplication between the feature vector $\mathbf{h}(\mathbf{x})$ and corresponding weight vector \mathbf{w} . We contrast the two outputs with and without sparsity:

$$f_c = \sum_{i=1}^m v_i \quad (\text{w.o sparsity}),$$

$$f_c^{\text{DICE}} = \sum_{i \in \text{top units}} v_i \quad (\text{w. sparsity}),$$

where f_c is the output using the summation of all units' contribution, and f_c^{DICE} takes the input from the top units (ranked based on the average contribution on ID data, see bottom of Figure 4).

DICE reduces the output variance We consider the unit contribution vector for OOD data $\mathbf{v} \in \mathbb{R}^m$, where each element is a *random variable* v_i with mean $\mathbb{E}[v_i] = \mu_i$ and variance $\text{Var}[v_i] = \sigma_i^2$. For simplicity, we assume each component is independent, but our theory can be extended to correlated variables (see Remark 1). Importantly, indices in \mathbf{v} are sorted based on *the same order* of unit contribution on ID data. By using units on the rightmost side, we now show the key result that DICE reduces the output variance.

Proposition 1. *Let v_i and v_j be two independent random variables. Denote the summation $r = v_i + v_j$, we have $\mathbb{E}[r] = \mathbb{E}[v_i] + \mathbb{E}[v_j]$ and $\text{Var}[r] = \text{Var}[v_i] + \text{Var}[v_j]$.*

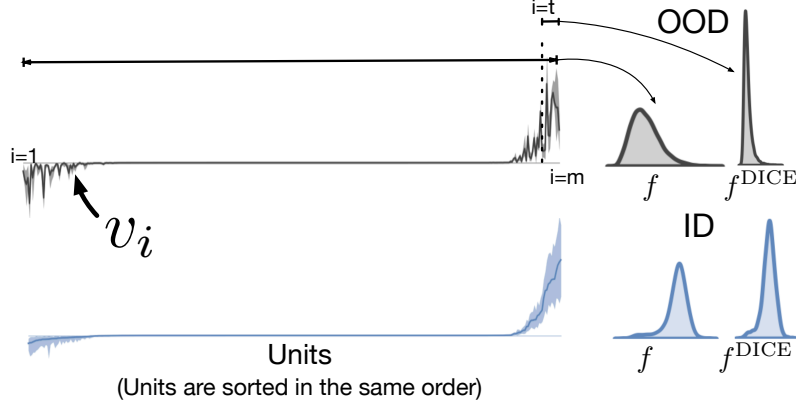


Fig. 4. Units in the penultimate layer are sorted based on the average contribution to a CIFAR-10 class (“airplane”). OOD data (SVHN) can trigger a non-negligible fraction of units with noisy signals on the CIFAR-10 trained model.

Lemma 1. *When taking the top $m - t$ units, the output variable f_c^{DICE} under sparsification has reduced variance:*

$$\text{Var}[f_c] - \text{Var}[f_c^{DICE}] = \sum_{i=1}^t \sigma_i^2$$

Proof. The proof directly follows Proposition 1.

Remark 1 (Extension to correlated variables) We can show in a more general case with correlated variables, the variance reduction is:

$$\sum_{i=1}^t \sigma_i^2 + 2 \sum_{1 \leq i < j \leq m} \text{Cov}(v_i, v_j) - 2 \sum_{t < i < j \leq m} \text{Cov}(v_i, v_j),$$

where $\text{Cov}(\cdot, \cdot)$ is the covariance. Our analysis shows that the covariance matrix primarily consists of 0, which indicates the independence of variables. Moreover, the summation of non-zero entries in the full matrix (i.e., the second term) is greater than that of the submatrix with top units (i.e., the third term), resulting in a larger variance reduction than in Lemma 1. See complete proof in Appendix E.

Remark 2 Energy score is compatible with DICE since it directly operates in the logit space. Our theoretical analysis above shows that DICE reduces the variance of each logit $f_c(\mathbf{x})$. This means that for detection scores such as energy score, the gap between OOD and ID score will be enlarged after applying DICE, which makes thresholding more capable of separating OOD and ID inputs and benefit OOD detection.

Table 5. Difference between the mean of ID’s output and OOD’s output. Here we use CIFAR-100 as ID data and $\Delta = \mathbb{E}_{\text{in}}[\max_c f_c^{\text{DICE}}] - \mathbb{E}_{\text{out}}[\max_c f_c^{\text{DICE}}]$ is averaged over six common OOD benchmark datasets described in Section 4

Sparsity	$p = 0.9$	$p = 0.7$	$p = 0.5$	$p = 0.3$	$p = 0.1$	$p = 0$
Δ	7.92	7.28	7.99	8.04	7.36	6.67

Remark 3 (Mean of output) Beyond variance, we further show in Table 5 the effect of sparsity on the mean of output: $\mathbb{E}_{\text{in}}[\max_c f_c^{\text{DICE}}]$ and $\mathbb{E}_{\text{out}}[\max_c f_c^{\text{DICE}}]$. The gap between the two directly translates into the OOD score separability. We show that DICE maintains similar (or even enlarges) differences in terms of mean as sparsity p increases. Therefore, DICE overall benefits OOD detection due to both *reduced output variances* and *increased differences of mean*—the combination of both effects leads to stronger separability between ID and OOD.

Remark 4 (Variance reduction on ID data) Note that we can also show the effect of variance reduction for ID data in a similar way. Importantly, DICE effectively preserves the most important information akin to the ID data, while reducing noisy signals that are harmful to OOD detection. Overall the variance reduction effect on both ID and OOD data leads to stronger separability.

7 Related Work

Out-of-distribution detection has attracted growing research attention in recent years. We highlight two major lines of work:

(1) One line of work perform OOD detection by devising scoring functions, including confidence-based methods [3, 20, 25, 34], energy-based score [35, 36, 44, 59, 51], distance-based approaches [32, 49, 52, 55], gradient-based score [24], and Bayesian approaches [11, 30, 38, 39, 40]. However, none of the previous methods considered weight sparsification for OOD detection. The closest work to ours is ReAct [51], which proposed truncating the high activations during test time for OOD detection. While ReAct only considers activation space, DICE examines both the weights and activation values together—the multiplication of which directly determines the unit contributions to the output. Our work is also related to [7], which pointed out that modern OOD detection methods succeed by detecting the existence of familiar features. DICE strengthens the familiarity hypothesis by keeping the dominating weights corresponding to the “major features”.

(2) A separate line of methods addressed OOD detection by training-time regularization [4, 5, 12, 18, 21, 26, 27, 31, 36, 39, 41, 42, 43, 56, 61, 65]. For example, models are encouraged to give predictions with uniform distribution [21, 31] or higher energies [8, 9, 27, 36, 42] for outlier data. The scope of this paper focuses on post hoc methods, which have the advantages of being easy to use and general applicability without modifying the training objective. The latter property is especially desirable for the adoption of OOD detection methods in real-world

production environments, when the overhead cost of retraining can be prohibitive.

Pruning and sparsification A great number of effort has been put into improving *post hoc* pruning and training time regularization for deep neural networks [12,13,15,16,33,37]. Many works obtain a sparse model by training with sparse regularization [12,15,37,53] or architecture modification [13,33], while our work primarily considers *post hoc* sparsification strategy which operates conveniently on a pre-trained network. On this line, two popular Bernoulli dropout techniques include unit dropout and weight dropout [50]. *Post hoc* pruning strategies truncate weights with low magnitude [16], or drop units with low weight norms [33]. In [62], they use a sparse linear layer to help identify spurious correlations and explain misclassifications. Orthogonal to existing works, our goal is to improve the OOD detection performance rather than accelerate computation and network debugging. In this paper, we first demonstrate that sparsification can be useful for OOD detection. An in-depth discussion and comparison of these methods are presented in Section 5.

Distributional shifts. Distributional shifts have attracted increasing research interest. It is important to recognize and differentiate various types of distributional shift problems. Literature in OOD detection is commonly concerned about model reliability and detection of semantic shifts, where the OOD inputs have disjoint labels *w.r.t.* ID data and therefore should not be predicted by the model. This is different from the OOD generalization task whose goal is to provide accurate predictions on OOD images under the same label space. For example, some works considered covariate shifts in the input space [28,19,47,54,69], where the model is expected to generalize to the OOD data.

8 Conclusion

This paper provides a simple sparsification strategy termed DICE, which ranks weights based on a contribution measure and then uses the most significant weights to derive the output for OOD detection. We provide both empirical and theoretical insights characterizing and explaining the mechanism by which DICE improves OOD detection. By exploiting the most important weights, DICE provably reduces the output variance for OOD data, resulting in a sharper output distribution and stronger separability from ID data. Extensive experiments show DICE can significantly improve the performance of OOD detection for over-parameterized networks. We hope our research can raise more attention to the importance of weight sparsification for OOD detection.

Acknowledgement

Work was supported by funding from Wisconsin Alumni Research Foundation (WARF). The authors would also like to thank reviewers for the helpful feedback.

References

1. Ba, J., Frey, B.: Adaptive dropout for training deep neural networks. In: Advances in Neural Information Processing Systems. vol. 26 (2013)
2. Babaeizadeh, M., Smaragdis, P., Campbell, R.H.: Noiseout: A simple way to prune neural networks. CoRR **abs/1611.06211** (2016)
3. Bendale, A., Boulton, T.E.: Towards open set deep networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1563–1572 (2016)
4. Bevanđić, P., Krešo, I., Oršić, M., Šegvić, S.: Discriminative out-of-distribution detection for semantic segmentation. arXiv preprint arXiv:1808.07703 (2018)
5. Chen, J., Li, Y., Wu, X., Liang, Y., Jha, S.: Atom: Robustifying out-of-distribution detection using outlier mining. In: Proceedings of European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (2021)
6. Cimpoi, M., Maji, S., Kokkinos, I., Mohamed, S., Vedaldi, A.: Describing textures in the wild. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 3606–3613 (2014)
7. Dietterich, T.G., Guyer, A.: The familiarity hypothesis: Explaining the behavior of deep open set methods. arXiv preprint arXiv:2203.02486 (2022)
8. Du, X., Wang, X., Gozum, G., Li, Y.: Unknown-aware object detection: Learning what you don’t know from videos in the wild. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2022)
9. Du, X., Wang, Z., Cai, M., Li, Y.: Vos: Learning what you don’t know by virtual outlier synthesis. In: Proceedings of the International Conference on Learning Representations (2022)
10. Filos, A., Tigkas, P., McAllister, R., Rhinehart, N., Levine, S., Gal, Y.: Can autonomous vehicles identify, recover from, and adapt to distribution shifts? In: Proceedings of the International Conference on Machine Learning. pp. 3145–3153. PMLR (2020)
11. Gal, Y., Ghahramani, Z.: Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In: Proceedings of the International Conference on Machine Learning. pp. 1050–1059 (2016)
12. Geifman, Y., El-Yaniv, R.: Selectivenet: A deep neural network with an integrated reject option. arXiv preprint arXiv:1901.09192 (2019)
13. Gomez, A.N., Zhang, I., Kamalakara, S.R., Madaan, D., Swersky, K., Gal, Y., Hinton, G.E.: Learning sparse networks using targeted dropout. arXiv preprint arXiv:1905.13678 (2019)
14. Goodfellow, I.J., Shlens, J., Szegedy, C.: Explaining and harnessing adversarial examples. arXiv preprint arXiv:1412.6572 (2014)
15. Han, S., Mao, H., Dally, W.J.: Deep compression: Compressing deep neural network with pruning, trained quantization and Huffman coding. In: Proceedings of the International Conference on Learning Representations (2016)
16. Han, S., Pool, J., Tran, J., Dally, W.: Learning both weights and connections for efficient neural network. In: Proceedings of the Advances in Neural Information Processing Systems. vol. 28, pp. 1135–1143 (2015)
17. He, K., Zhang, X., Ren, S., Sun, J.: Identity mappings in deep residual networks. In: Proceedings of the European conference on computer vision. pp. 630–645. Springer (2016)
18. Hein, M., Andriushchenko, M., Bitterwolf, J.: Why relu networks yield high-confidence predictions far away from the training data and how to mitigate the problem. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 41–50 (2019)

19. Hendrycks, D., Dietterich, T.: Benchmarking neural network robustness to common corruptions and perturbations. arXiv preprint arXiv:1903.12261 (2019)
20. Hendrycks, D., Gimpel, K.: A baseline for detecting misclassified and out-of-distribution examples in neural networks. *Proceedings of International Conference on Learning Representations* (2017)
21. Hendrycks, D., Mazeika, M., Dietterich, T.: Deep anomaly detection with outlier exposure. arXiv preprint arXiv:1812.04606 (2018)
22. Hsu, Y.C., Shen, Y., Jin, H., Kira, Z.: Generalized odin: Detecting out-of-distribution image without learning from out-of-distribution data. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2020)
23. Huang, G., Liu, Z., Van Der Maaten, L., Weinberger, K.Q.: Densely connected convolutional networks. In: *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*. pp. 4700–4708 (2017)
24. Huang, R., Geng, A., Li, Y.: On the importance of gradients for detecting distributional shifts in the wild. In: *Proceedings of the Advances in Neural Information Processing Systems* (2021)
25. Huang, R., Li, Y.: Towards scaling out-of-distribution detection for large semantic space. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2021)
26. Jeong, T., Kim, H.: Ood-maml: Meta-learning for few-shot out-of-distribution detection and classification. In: *Proceedings of the Advances in Neural Information Processing Systems* (2020)
27. Katz-Samuels, J., Nakhleh, J., Nowak, R., Li, Y.: Training ood detectors in their natural habitats. In: *Proceedings of the International Conference on Machine Learning*. PMLR (2022)
28. Koh, P.W., Sagawa, S., Xie, S.M., Zhang, M., Balsubramani, A., Hu, W., Yasunaga, M., Phillips, R.L., Gao, I., Lee, T., et al.: Wilds: A benchmark of in-the-wild distribution shifts. In: *Proceedings of the International Conference on Machine Learning*. pp. 5637–5664. PMLR (2021)
29. Krizhevsky, A., Hinton, G., et al.: Learning multiple layers of features from tiny images (2009)
30. Lakshminarayanan, B., Pritzel, A., Blundell, C.: Simple and scalable predictive uncertainty estimation using deep ensembles. In: *Advances in neural information processing systems*. pp. 6402–6413 (2017)
31. Lee, K., Lee, H., Lee, K., Shin, J.: Training confidence-calibrated classifiers for detecting out-of-distribution samples. arXiv preprint arXiv:1711.09325 (2017)
32. Lee, K., Lee, K., Lee, H., Shin, J.: A simple unified framework for detecting out-of-distribution samples and adversarial attacks. In: *Advances in Neural Information Processing Systems*. pp. 7167–7177 (2018)
33. Li, H., Kadav, A., Durdanovic, I., Samet, H., Graf, H.P.: Pruning filters for efficient convnets. In: *Proceedings of International Conference on Learning Representations* (2017)
34. Liang, S., Li, Y., Srikant, R.: Enhancing the reliability of out-of-distribution image detection in neural networks. In: *Proceedings of International Conference on Learning Representations* (2018)
35. Lin, Z., Roy, S.D., Li, Y.: Mood: Multi-level out-of-distribution detection. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 15313–15323 (June 2021)
36. Liu, W., Wang, X., Owens, J., Li, Y.: Energy-based out-of-distribution detection. In: *Proceedings of the Advances in Neural Information Processing Systems* (2020)

37. Louizos, C., Welling, M., Kingma, D.P.: Learning sparse neural networks through l_0 regularization. In: International Conference on Learning Representations (2018)
38. Maddox, W.J., Izmailov, P., Garipov, T., Vetrov, D.P., Wilson, A.G.: A simple baseline for bayesian uncertainty in deep learning. *Advances in Neural Information Processing Systems* **32**, 13153–13164 (2019)
39. Malinin, A., Gales, M.: Predictive uncertainty estimation via prior networks. In: *Advances in Neural Information Processing Systems*. pp. 7047–7058 (2018)
40. Malinin, A., Gales, M.: Reverse kl-divergence training of prior networks: Improved uncertainty and adversarial robustness. In: *Advances in Neural Information Processing Systems* (2019)
41. Meinke, A., Hein, M.: Towards neural networks that provably know when they don’t know. *arXiv preprint arXiv:1909.12180* (2019)
42. Ming, Y., Fan, Y., Li, Y.: Poem: Out-of-distribution detection with posterior sampling. In: *Proceedings of the International Conference on Machine Learning*. PMLR (2022)
43. Mohseni, S., Pitale, M., Yadawa, J., Wang, Z.: Self-supervised learning for generalizable out-of-distribution detection. In: *AAAI*. pp. 5216–5223 (2020)
44. Morteza, P., Li, Y.: Provable guarantees for understanding out-of-distribution detection. *Proceedings of the AAAI Conference on Artificial Intelligence* (2022)
45. Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B., Ng, A.Y.: Reading digits in natural images with unsupervised feature learning (2011)
46. Nguyen, A., Yosinski, J., Clune, J.: Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 427–436 (2015)
47. Ovadia, Y., Fertig, E., Ren, J., Nado, Z., Sculley, D., Nowozin, S., Dillon, J., Lakshminarayanan, B., Snoek, J.: Can you trust your model’s uncertainty? evaluating predictive uncertainty under dataset shift. In: *Proceedings of the Advances in Neural Information Processing Systems*. vol. 32, pp. 13991–14002 (2019)
48. Roy, A.G., Ren, J., Azizi, S., Loh, A., Natarajan, V., Mustafa, B., Pawlowski, N., Freyberg, J., Liu, Y., Beaver, Z., et al.: Does your dermatology classifier know what it doesn’t know? detecting the long-tail of unseen conditions. *arXiv preprint arXiv:2104.03829* (2021)
49. Schwag, V., Chiang, M., Mittal, P.: Ssd: A unified framework for self-supervised outlier detection. In: *International Conference on Learning Representations* (2021)
50. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., Salakhutdinov, R.: Dropout: A simple way to prevent neural networks from overfitting. In: *Journal of Machine Learning Research*. vol. 15, pp. 1929–1958 (2014)
51. Sun, Y., Guo, C., Li, Y.: React: Out-of-distribution detection with rectified activations. In: *Advances in Neural Information Processing Systems* (2021)
52. Sun, Y., Ming, Y., Zhu, X., Li, Y.: Out-of-distribution detection with deep nearest neighbors. In: *Proceedings of the International Conference on Machine Learning* (2022)
53. Sun, Y., Ravi, S., Singh, V.: Adaptive activation thresholding: Dynamic routing type behavior for interpretability in convolutional neural networks. In: *Proceedings of the International Conference on Computer Vision* (2019)
54. Sun, Y., Wang, X., Liu, Z., Miller, J., Efros, A., Hardt, M.: Test-time training with self-supervision for generalization under distribution shifts. In: *Proceedings of the International Conference on Machine Learning*. pp. 9229–9248. PMLR (2020)
55. Tack, J., Mo, S., Jeong, J., Shin, J.: Csi: Novelty detection via contrastive learning on distributionally shifted instances. In: *Advances in Neural Information Processing Systems* (2020)

56. Van Amersfoort, J., Smith, L., Teh, Y.W., Gal, Y.: Uncertainty estimation using a single deep deterministic neural network. In: *Proceedings of the International Conference on Machine Learning* (2020)
57. Van Horn, G., Mac Aodha, O., Song, Y., Cui, Y., Sun, C., Shepard, A., Adam, H., Perona, P., Belongie, S.: The inaturalist species classification and detection dataset. In: *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*. pp. 8769–8778 (2018)
58. Wan, L., Zeiler, M.D., Zhang, S., LeCun, Y., Fergus, R.: Regularization of neural networks using dropconnect. In: *Proceedings of the International Conference on Machine Learning*. vol. 28, pp. 1058–1066 (2013)
59. Wang, H., Liu, W., Bocchieri, A., Li, Y.: Can multi-label classification networks know what they don’t know? *Proceedings of the Advances in Neural Information Processing Systems* (2021)
60. Wang, X., Peng, Y., Lu, L., Lu, Z., Bagheri, M., Summers, R.M.: Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 2097–2106 (2017)
61. Wei, H., Xie, R., Cheng, H., Feng, L., An, B., Li, Y.: Mitigating neural network overconfidence with logit normalization. *Proceedings of the International Conference on Machine Learning* (2022)
62. Wong, E., Santurkar, S., Madry, A.: Leveraging sparse linear layers for debuggable deep networks. In: *Proceedings of the International Conference on Machine Learning*. pp. 11205–11216. PMLR (2021)
63. Xiao, J., Hays, J., Ehinger, K.A., Oliva, A., Torralba, A.: Sun database: Large-scale scene recognition from abbey to zoo. In: *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*. pp. 3485–3492. IEEE Computer Society (2010)
64. Xu, P., Ehinger, K.A., Zhang, Y., Finkelstein, A., Kulkarni, S.R., Xiao, J.: Turkergaze: Crowdsourcing saliency with webcam based eye tracking. *arXiv preprint arXiv:1504.06755* (2015)
65. Yang, J., Wang, H., Feng, L., Yan, X., Zheng, H., Zhang, W., Liu, Z.: Semantically coherent out-of-distribution detection. In: *Proceedings of the IEEE International Conference on Computer Vision*. pp. 8301–8309 (October 2021)
66. Yu, F., Seff, A., Zhang, Y., Song, S., Funkhouser, T., Xiao, J.: Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv preprint arXiv:1506.03365* (2015)
67. Zhang, C., Bengio, S., Hardt, M., Recht, B., Vinyals, O.: Understanding deep learning requires rethinking generalization. In: *Proceedings of International Conference on Learning Representations*
68. Zhou, B., Lapedriza, A., Khosla, A., Oliva, A., Torralba, A.: Places: A 10 million image database for scene recognition. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence*. vol. 40, pp. 1452–1464. IEEE (2017)
69. Zhou, K., Liu, Z., Qiao, Y., Xiang, T., Loy, C.C.: Domain generalization: A survey (2021)