

CSI: Novelty Detection via Contrastive Learning on Distributionally Shifted Instances

Jihoon Tack^{*†}, Sangwoo Mo^{*‡}, Jongheon Jeong[‡], Jinwoo Shin^{†‡}

[†]Graduate School of AI, KAIST

[‡]School of Electrical Engineering, KAIST

{jihoontack, swmo, jongheonj, jinwoos}@kaist.ac.kr

Abstract

Novelty detection, *i.e.*, identifying whether a given sample is drawn from outside the training distribution, is essential for reliable machine learning. To this end, there have been many attempts at learning a representation well-suited for novelty detection and designing a score based on such representation. In this paper, we propose a simple, yet effective method named *contrasting shifted instances* (CSI), inspired by the recent success on contrastive learning of visual representations. Specifically, in addition to contrasting a given sample with other instances as in conventional contrastive learning methods, our training scheme contrasts the sample with distributionally-shifted augmentations of itself. Based on this, we propose a new detection score that is specific to the proposed training scheme. Our experiments demonstrate the superiority of our method under various novelty detection scenarios, including unlabeled one-class, unlabeled multi-class and labeled multi-class settings, with various image benchmark datasets. Code and pre-trained models are available at <https://github.com/alinlab/CSI>.

1 Introduction

Out-of-distribution (OOD) detection [26], also referred to as a novelty- or anomaly detection is a task of identifying whether a test input is drawn far from the training distribution (in-distribution) or not. In general, the OOD detection problem aims to detect OOD samples where a detector is allowed to access only to training data. The space of OOD samples is typically huge, *i.e.*, an OOD sample can vary significantly and arbitrarily from the given training distribution. Hence, assuming specific prior knowledge, *e.g.*, external data representing some specific OODs, may introduce a bias to the detector. The OOD detection is a classic yet essential problem in machine learning, with a broad range of applications, including medical diagnosis [4], fraud detection [53], and autonomous driving [12].

A long line of literature has thus been proposed, including density-based [74, 46, 6, 47, 11, 55, 61, 17], reconstruction-based [58, 76, 9, 54, 52, 7], one-class classifier [59, 56], and self-supervised [15, 25, 2] approaches. Overall, a majority of recent literature is concerned with (a) modeling the representation to better encode normality [23, 25], and (b) defining a new detection score [56, 2]. In particular, recent studies have shown that inductive biases from self-supervised learning significantly help to learn discriminative features for OOD detection [15, 25, 2].

Meanwhile, recent progress on self-supervised learning has proven the effectiveness of *contrastive learning* in various domains, *e.g.*, computer vision [21, 5], audio processing [50], and reinforcement learning [63]. Contrastive learning extracts a strong inductive bias from multiple (similar) views of a sample by let them attract each other, yet repelling them to other samples. *Instance discrimination* [69]

^{*}Equal contribution

is a special type of contrastive learning where the views are restricted up to different augmentations, which have achieved state-of-the-art results on visual representation learning [21, 5].

Inspired by the recent success of instance discrimination, we aim to utilize its power of representation learning for OOD detection. To this end, we investigate the following questions: (a) how to learn a (more) discriminative representation for detecting OODs and (b) how to design a score function utilizing the representation from (a). We remark that the desired representation for OOD detection may differ from that for standard representation learning [23, 25], as the former aims to discriminate in-distribution and OOD samples, while the latter aims to discriminate *within* in-distribution samples.

We first found that existing contrastive learning scheme is already reasonably effective for detecting OOD samples with a proper detection score. We further observe that one can improve its performance by utilizing “hard” augmentations, *e.g.*, rotation, that were known to be harmful and unused for the standard contrastive learning [5]. In particular, while the existing contrastive learning schemes act by pulling all augmented samples toward the original sample, we suggest to additionally push the samples with hard or distribution-shifting augmentations away from the original. We observe that contrasting shifted samples help OOD detection, as the model now learns a new task of discriminating *between* in- and out-of-distribution, in addition to the original task of discriminating *within* in-distribution.

Contribution. We propose a simple yet effective method for OOD detection, coined *contrasting shifted instances* (CSI). Built upon the existing contrastive learning scheme [5], we propose two novel additional components: (a) a new training method which contrasts distributionally-shifted augmentations (of the given sample) in addition to other instances, and (b) a score function which utilizes both the contrastively learned representation and our new training scheme in (a). Finally, we show that CSI enjoys broader usage by applying it to improve the confidence-calibration of the classifiers: it relaxes the overconfidence issue in their predictions for both in- and out-of-distribution samples while maintaining the classification accuracy.

We verify the effectiveness of CSI under various environments of detecting OOD, including unlabeled one-class, unlabeled multi-class, and labeled multi-class settings. To our best knowledge, we are the first to demonstrate all three settings under a single framework. Overall, CSI outperforms the baseline methods for all tested datasets. In particular, CSI achieves new state-of-the-art results² on one-class classification, *e.g.*, it improves the mean area under the receiver operating characteristics (AUROC) from 90.1% to 94.3% (+4.2%) for CIFAR-10 [33], 79.8% to 89.6% (+9.8%) for CIFAR-100 [33], and 85.7% to 91.6% (+5.9%) for ImageNet-30 [25] one-class datasets, respectively. We remark that CSI gives a larger improvement in harder (or near-distribution) OOD samples. To verify this, we also release new benchmark datasets: fixed version of the resized LSUN and ImageNet [39].

We remark that learning representation to discriminate in- vs. out-of-distributions is an important but under-explored problem. We believe that our work would guide new interesting directions in the future, for both representation learning and OOD detection.

2 CSI: Contrasting shifted instances

For a given dataset $\{x_m\}_{m=1}^M$ sampled from a data distribution $p_{\text{data}}(x)$ on the data space \mathcal{X} , the goal of out-of-distribution (OOD) detection is to model a detector from $\{x_m\}$ that identifies whether x is sampled from the data generating distribution (or in-distribution) $p_{\text{data}}(x)$ or not. As modeling $p_{\text{data}}(x)$ directly is prohibitive in most cases, many existing methods for OOD detection define a *score function* $s(x)$ that a high value heuristically represents that x is from in-distribution.

2.1 Contrastive learning

The idea of *contrastive learning* is to learn an encoder f_θ to extract the necessary information to distinguish similar samples from the others. Let x be a query, $\{x_+\}$, and $\{x_-\}$ be a set of positive and negative samples, respectively, and $\text{sim}(z, z') := z \cdot z' / \|z\| \|z'\|$ be the cosine similarity. Then, the primitive form of the *contrastive loss* is defined as follows:

$$\mathcal{L}_{\text{con}}(x, \{x_+\}, \{x_-\}) := -\frac{1}{|\{x_+\}|} \log \frac{\sum_{x' \in \{x_+\}} \exp(\text{sim}(z(x), z(x'))/\tau)}{\sum_{x' \in \{x_+\} \cup \{x_-\}} \exp(\text{sim}(z(x), z(x'))/\tau)}, \quad (1)$$

where $|\{x_+\}|$ denotes the cardinality of the set $\{x_+\}$, $z(x)$ denotes the output feature of the contrastive layer, and τ denotes a temperature hyper-parameter. One can define the contrastive feature $z(x)$

²We do not compare with methods using *external* OOD samples [24, 57].

directly from the encoder f_θ , *i.e.*, $z(x) = f_\theta(x)$ [21], or apply an additional projection layer g_ϕ , *i.e.*, $z(x) = g_\phi(f_\theta(x))$ [5]. We use the projection layer following the recent studies [5, 30].

In this paper, we specifically consider the simple contrastive learning (*SimCLR*) [5], a simple and effective objective based on the task of *instance discrimination* [69]. Let $\tilde{x}_i^{(1)}$ and $\tilde{x}_i^{(2)}$ be two independent augmentations of x_i from a pre-defined family \mathcal{T} , namely, $\tilde{x}_i^{(1)} := T_1(x_i)$ and $\tilde{x}_i^{(2)} := T_2(x_i)$, where $T_1, T_2 \sim \mathcal{T}$. Then the SimCLR objective can be defined by the contrastive loss (1) where each $(\tilde{x}_i^{(1)}, \tilde{x}_i^{(2)})$ and $(\tilde{x}_i^{(2)}, \tilde{x}_i^{(1)})$ are considered as query-key pairs while others being negatives. Namely, for a given batch $\mathcal{B} := \{x_i\}_{i=1}^B$, the SimCLR objective is defined as follows:

$$\mathcal{L}_{\text{SimCLR}}(\mathcal{B}; \mathcal{T}) := \frac{1}{2B} \sum_{i=1}^B \mathcal{L}_{\text{con}}(\tilde{x}_i^{(1)}, \tilde{x}_i^{(2)}, \tilde{\mathcal{B}}_{-i}) + \mathcal{L}_{\text{con}}(\tilde{x}_i^{(2)}, \tilde{x}_i^{(1)}, \tilde{\mathcal{B}}_{-i}), \quad (2)$$

where $\tilde{\mathcal{B}} := \{\tilde{x}_i^{(1)}\}_{i=1}^B \cup \{\tilde{x}_i^{(2)}\}_{i=1}^B$ and $\tilde{\mathcal{B}}_{-i} := \{\tilde{x}_j^{(1)}\}_{j \neq i} \cup \{\tilde{x}_j^{(2)}\}_{j \neq i}$.

2.2 Contrastive learning for distribution-shifting transformations

Chen et al. [5] has performed an extensive study on which family of augmentations \mathcal{T} leads to a better representation when used in SimCLR, *i.e.*, which transformations should f_θ consider as positives. Overall, the authors report that some of the examined augmentations (*e.g.*, rotation), sometimes degrades the discriminative performance of SimCLR. One of our key findings is that such augmentations can be useful for OOD detection by considering them as *negatives* - contrast from the original sample. In this paper, we explore which family of augmentations \mathcal{S} , which we call *distribution-shifting transformations*, or simply *shifting transformations*, would lead to better representation in terms of OOD detection when used as negatives in SimCLR.

Contrasting shifted instances. We consider a set \mathcal{S} consisting of K different (random or deterministic) transformations, including the identity I : namely, we denote $\mathcal{S} := \{S_0 = I, S_1, \dots, S_{K-1}\}$. In contrast to the vanilla SimCLR that considers augmented samples as positive to each other, we attempt to consider them as negative if the augmentation is from \mathcal{S} . For a given batch of samples $\mathcal{B} = \{x_i\}_{i=1}^B$, this can be done simply by augmenting \mathcal{B} via \mathcal{S} before putting it into the SimCLR loss defined in (2): namely, we define *contrasting shifted instances* (con-SI) loss as follows:

$$\mathcal{L}_{\text{con-SI}} := \mathcal{L}_{\text{SimCLR}}\left(\bigcup_{S \in \mathcal{S}} \mathcal{B}_S; \mathcal{T}\right), \quad \text{where } \mathcal{B}_S := \{S(x_i)\}_{i=1}^B. \quad (3)$$

Here, our intuition is to regard each distributionally-shifted sample (*i.e.*, $S \neq I$) as an OOD with respect to the original. In this respect, con-SI attempts to discriminate an in-distribution (*i.e.*, $S = I$) sample from other OOD (*i.e.*, $S \in \{S_1, \dots, S_{K-1}\}$) samples. We further verify the effectiveness of con-SI in our experimental results: although con-SI does not improve representation for standard classification, it does improve OOD detection significantly (see linear evaluation in Section 3.2).

Classifying shifted instances. In addition to contrasting shifted instances, we consider an auxiliary task that predicts which shifting transformation $y^S \in \mathcal{S}$ is applied for a given input x , in order to facilitate f_θ to discriminate each shifted instance. Specifically, we add a linear layer to f_θ for modeling an auxiliary softmax classifier $p_{\text{cls-SI}}(y^S | x)$, as in [15, 25, 2]. Let $\tilde{\mathcal{B}}_S$ be the batch augmented from \mathcal{B}_S via SimCLR; then, we define *classifying shifted instances* (cls-SI) loss as follows:

$$\mathcal{L}_{\text{cls-SI}} := \frac{1}{2B} \frac{1}{K} \sum_{S \in \mathcal{S}} \sum_{\tilde{x}_S \in \tilde{\mathcal{B}}_S} -\log p_{\text{cls-SI}}(y^S = S | \tilde{x}_S). \quad (4)$$

The final loss of our proposed method, *CSI*, is defined by combining the two objectives:

$$\mathcal{L}_{\text{CSI}} = \mathcal{L}_{\text{con-SI}} + \lambda \cdot \mathcal{L}_{\text{cls-SI}} \quad (5)$$

where $\lambda > 0$ is a balancing hyper-parameter. We simply set $\lambda = 1$ for all our experiments.

OOD-ness: How to choose the shifting transformation? In principle, we choose the shifting transformation that generates the most OOD-like yet semantically meaningful samples. Intuitively, such samples can be most effective (‘nearby’ but ‘not-too-nearby’) OOD samples, as also discussed in Section 3.2. More specifically, we measure the *OOD-ness* of a transformation by the area under the receiver operating characteristics (AUROC) between in-distribution vs. transformed samples under vanilla SimCLR, using the detection score (6) defined in Section 2.3. The transformation with high OOD-ness values (*i.e.*, OOD-like) indeed performs better (see Table 4 and Table 5 in Section 3.2).

2.3 Score functions for detecting out-of-distribution

Upon the representation $z(\cdot)$ learned by our proposed training objective, we define several score functions for detecting out-of-distribution; whether a given x is OOD or not. We first propose a detection score that is applicable to any contrastive representation. We then introduce how one could incorporate additional information learned by contrasting (and classifying) shifted instances as in (5).

Detection score for contrastive representation. Overall, we find that two features from SimCLR representations are surprisingly effective for detecting OOD samples: (a) the *cosine similarity* to the nearest training sample in $\{x_m\}$, i.e., $\max_m \text{sim}(z(x_m), z(x))$, and (b) the *norm* of the representation, i.e., $\|z(x)\|$. Intuitively, the contrastive loss increases the norm of in-distribution samples, as it is an easy way to minimize the cosine similarity of identical samples by increasing the denominator of (1). We discuss further detailed analysis of both features in Appendix H. We simply combine these features to define a detection score s_{con} for contrastive representation:

$$s_{\text{con}}(x; \{x_m\}) := \max_m \text{sim}(z(x_m), z(x)) \cdot \|z(x)\|. \quad (6)$$

We also discuss how one can reduce the computation and memory cost by choosing a proper subset (i.e., coreset) of training samples in Appendix E.

Utilizing shifting transformations. Given that our proposed \mathcal{L}_{CSI} is used for training, one can further improve the detection score s_{con} significantly by incorporating shifting transformations \mathcal{S} . Here, we propose two additional scores, $s_{\text{con-SI}}$ and $s_{\text{cls-SI}}$, where are corresponded to $\mathcal{L}_{\text{con-SI}}$ (3) and $\mathcal{L}_{\text{cls-SI}}$ (4), respectively.

Firstly, we define $s_{\text{con-SI}}$ by taking an expectation of s_{con} over $S \in \mathcal{S}$:

$$s_{\text{con-SI}}(x; \{x_m\}) := \sum_{S \in \mathcal{S}} \lambda_S^{\text{con}} s_{\text{con}}(S(x); \{S(x_m)\}), \quad (7)$$

where $\lambda_S^{\text{con}} := M / \sum_m s_{\text{con}}(S(x_m); \{S(x_m)\}) = M / \sum_m \|z(S(x_m))\|$ for M training samples is a balancing term to scale the scores of each shifting transformation (See Appendix F for details).

Secondly, we define $s_{\text{cls-SI}}$ utilizing the auxiliary classifier $p(y^S|x)$ upon f_θ as follows:

$$s_{\text{cls-SI}}(x) := \sum_{S \in \mathcal{S}} \lambda_S^{\text{cls}} W_S f_\theta(S(x)), \quad (8)$$

where $\lambda_S^{\text{cls}} := M / \sum_m [W_S f_\theta(S(x_m))]$ are again balancing terms similarly to above, and W_S is the weight vector in the linear layer of $p(y^S|x)$ per $S \in \mathcal{S}$.

Finally, the combined score for CSI representation is defined as follows:

$$s_{\text{CSI}}(x; \{x_m\}) := s_{\text{con-SI}}(x; \{x_m\}) + s_{\text{cls-SI}}(x). \quad (9)$$

Ensembling over random augmentations. In addition, we find one can further improve each of the proposed scores by ensembling it over random augmentations $T(x)$ where $T \sim \mathcal{T}$. Namely, for instance, the *ensembled* CSI score is defined by $s_{\text{CSI-ens}}(x) := \mathbb{E}_{T \sim \mathcal{T}}[s_{\text{CSI}}(T(x))]$. Unless otherwise noted, we use these ensembled versions of (6) to (9) in our experiments. See Appendix D for details.

2.4 Extension for training confidence-calibrated classifiers

Furthermore, we propose an extension of CSI for training *confidence-calibrated* classifiers [22, 37] from a given labeled dataset $\{(x_m, y_m)\}_m \subseteq \mathcal{X} \times \mathcal{Y}$ by adapting it to *supervised contrastive learning* (SupCLR) [30]. Here, the goal is to model a classifier $p(y|x)$ that is (a) accurate on predicting y when x is in-distribution, and (b) the *confidence* $s_{\text{sup}}(x) := \max_y p(y|x)$ [22] of the classifier is *well-calibrated*, i.e., $s_{\text{sup}}(x)$ should be low if x is an OOD sample or $\arg \max_y p(y|x) \neq \text{true label}$.

Supervised contrastive learning (SupCLR). SupCLR is a supervised extension of SimCLR that contrasts samples in *class-wise*, instead of in instance-wise: every samples of the same classes are considered as positives. Let $\mathcal{C} = \{(x_i, y_i)\}_{i=1}^B$ be a training batch with class labels $y_i \in \mathcal{Y}$, and $\tilde{\mathcal{C}}$ be an augmented batch by random transformation \mathcal{T} , i.e., $\tilde{\mathcal{C}} := \{(\tilde{x}_j, y_j) \mid \tilde{x}_j \in \tilde{\mathcal{B}}\}$. For a given label y , we divide $\tilde{\mathcal{C}}$ into two subsets $\tilde{\mathcal{C}} = \tilde{\mathcal{C}}_y \cup \tilde{\mathcal{C}}_{-y}$ where $\tilde{\mathcal{C}}_y$ contains the samples of label y and $\tilde{\mathcal{C}}_{-y}$ contains the remaining. Then, the SupCLR objective is defined by:

$$\mathcal{L}_{\text{SupCLR}}(\mathcal{C}; \mathcal{T}) := \frac{1}{2B} \sum_{j=1}^{2B} \mathcal{L}_{\text{con}}(\tilde{x}_j, \tilde{\mathcal{C}}_{y_j} \setminus \{\tilde{x}_j\}, \tilde{\mathcal{C}}_{-y_j}). \quad (10)$$

Table 1: AUROC (%) of various OOD detection methods trained on one-class dataset of (a) CIFAR-10, (b) CIFAR-100 (super-class), and (c) ImageNet-30. For CIFAR-10, we report the means and standard deviations of per-class AUROC averaged over five trials, and the final column indicates the mean AUROC across all the classes. For CIFAR-100 and ImageNet-30, we only report the mean AUROC over a single trial. Bold denotes the best results, and * denotes the values from the reference. See Appendix C for additional results, *e.g.*, per-class AUROC on CIFAR-100 and ImageNet-30.

(a) One-class CIFAR-10												
Method	Network	Plane	Car	Bird	Cat	Deer	Dog	Frog	Horse	Ship	Truck	Mean
OC-SVM* [59]	-	65.6	40.9	65.3	50.1	75.2	51.2	71.8	51.2	67.9	48.5	58.8
DeepSVDD* [56]	LeNet	61.7	65.9	50.8	59.1	60.9	65.7	67.7	67.3	75.9	73.1	64.8
AnoGAN* [58]	DCGAN	67.1	54.7	52.9	54.5	65.1	60.3	58.5	62.5	75.8	66.5	61.8
OCGAN* [52]	OCGAN	75.7	53.1	64.0	62.0	72.3	62.0	72.3	57.5	82.0	55.4	65.7
Geom* [15]	WRN-16-8	74.7	95.7	78.1	72.4	87.8	87.8	83.4	95.5	93.3	91.3	86.0
Rot* [25]	WRN-16-4	71.9	94.5	78.4	70.0	77.2	86.6	81.6	93.7	90.7	88.8	83.3
Rot+Trans* [25]	WRN-16-4	77.5	96.9	87.3	80.9	92.7	90.2	90.9	96.5	95.2	93.3	90.1
GOAD* [2]	WRN-10-4	77.2	96.7	83.3	77.7	87.8	87.8	90.0	96.1	93.8	92.0	88.2
Rot [25]	ResNet-18	78.3±0.2	94.3±0.3	86.2±0.4	80.8±0.6	89.4±0.5	89.0±0.4	88.9±0.4	95.1±0.2	92.3±0.3	89.7±0.3	88.4
Rot+Trans [25]	ResNet-18	80.4±0.3	96.4±0.2	85.9±0.3	81.1±0.5	91.3±0.3	89.6±0.3	89.9±0.3	95.9±0.1	95.0±0.1	92.6±0.2	89.8
GOAD [2]	ResNet-18	75.5±0.3	94.1±0.3	81.8±0.5	72.0±0.3	83.7±0.9	84.4±0.3	82.9±0.8	93.9±0.3	92.9±0.3	89.5±0.2	85.1
CSI (ours)	ResNet-18	89.9±0.1	99.1±0.0	93.1±0.2	86.4±0.2	93.9±0.1	93.2±0.2	95.1±0.1	98.7±0.0	97.9±0.0	95.5±0.1	94.3
(b) One-class CIFAR-100 (super-class)						(c) One-class ImageNet-30						
Method	Network	AUROC				Method	Network	AUROC				
OC-SVM* [59]	-	63.1				Rot* [25]	ResNet-18	65.3				
Geom* [15]	WRN-16-8	78.7				Rot+Trans* [25]	ResNet-18	77.9				
Rot [25]	ResNet-18	77.7				Rot+Attn* [25]	ResNet-18	81.6				
Rot+Trans [25]	ResNet-18	79.8				Rot+Trans+Attn* [25]	ResNet-18	84.8				
GOAD [2]	ResNet-18	74.5				Rot+Trans+Attn+Resize* [25]	ResNet-18	85.7				
CSI (ours)	ResNet-18	89.6				CSI (ours)	ResNet-18	91.6				

After training the embedding network $f_\theta(x)$ with the SupCLR objective (10), we train a linear classifier upon $f_\theta(x)$ to model $p_{\text{SupCLR}}(y|x)$.

Supervised extension of CSI. We extend CSI by incorporating the shifting transformations \mathcal{S} into the SupCLR objective: here, we consider a joint label $(y, y^S) \in \mathcal{Y} \times \mathcal{S}$ of class label y and shifting transformation y^S . Then, the *supervised contrasting shifted instances* (sup-CSI) loss is given by:

$$\mathcal{L}_{\text{sup-CSI}} := \mathcal{L}_{\text{SupCLR}} \left(\bigcup_{S \in \mathcal{S}} \mathcal{C}_S; \mathcal{T} \right), \quad \text{where } \mathcal{C}_S := \{(S(x_i), (y_i, S))\}_{i=1}^B. \quad (11)$$

Note that we do not use the auxiliary classification loss $\mathcal{L}_{\text{cls-SI}}$ (4), since the objective already classifies the shifted instances under a *self-label augmented* [35] space $\mathcal{Y} \times \mathcal{S}$.

Upon the learned representation via (11), we additionally train two linear classifiers: $p_{\text{CSI}}(y|x)$ and $p_{\text{CSI-joint}}(y, y^S|x)$ that predicts the class labels and joint labels, respectively. We directly apply $s_{\text{sup}}(x)$ for the former $p_{\text{CSI}}(y|x)$. For the latter, on the other hand, we marginalize the joint prediction over the shifting transformation in a similar manner of Section 2.3. Precisely, let $l(x) \in \mathbb{R}^{C \times K}$ be logit values of $p_{\text{CSI-joint}}(y, y^S|x)$ for $|\mathcal{Y}| = C$ and $|\mathcal{S}| = K$, and $l(x)_k \in \mathbb{R}^C$ be logit values correspond to $p_{\text{CSI-joint}}(y, y^S = S_k|x)$. Then, the ensembled probability is:

$$p_{\text{CSI-ens}}(y|x) := \sigma \left(\frac{1}{K} \sum_k l(S_k(x))_k \right), \quad (12)$$

where σ denotes the softmax activation. Here, we use $p_{\text{CSI-ens}}$ to compute the confidence $s_{\text{sup}}(x)$. We denote the confidence computed by p_{CSI} and $p_{\text{CSI-ens}}$ and “CSI” and “CSI-ens”, respectively.

3 Experiments

In Section 3.1, we report OOD detection results on unlabeled one-class, unlabeled multi-class, and labeled multi-class datasets. In Section 3.2, we analyze the effects on various shifting transformations in the context of OOD detection, as well as an ablation study on each component we propose.

Table 2: AUROC (%) of various OOD detection methods trained on unlabeled (a) CIFAR-10 and (b) ImageNet-30. The reported results are averaged over five trials, subscripts denote standard deviation, and bold denote the best results. * denotes the values from the reference.

(a) Unlabeled CIFAR-10								
CIFAR10 →								
Method	Network	SVHN	LSUN	ImageNet	LSUN (FIX)	ImageNet (FIX)	CIFAR-100	Interp.
Likelihood*	PixelCNN++	8.3	-	64.2	-	-	52.6	52.6
Likelihood*	Glow	8.3	-	66.3	-	-	58.2	58.2
Likelihood*	EBM	63.0	-	-	-	-	-	70.0
Likelihood Ratio* [55]	PixelCNN++	91.2	-	-	-	-	-	-
Input Complexity* [61]	PixelCNN++	92.9	-	58.9	-	-	53.5	-
Input Complexity* [61]	Glow	95.0	-	71.6	-	-	73.6	-
Rot [25]	ResNet-18	97.6±0.2	89.2±0.7	90.5±0.3	77.7±0.3	83.2±0.1	79.0±0.1	64.0±0.3
Rot+Trans [25]	ResNet-18	97.8±0.2	92.8±0.9	94.2±0.7	81.6±0.4	86.7±0.1	82.3±0.2	68.1±0.8
GOAD [2]	ResNet-18	96.3±0.2	89.3±1.5	91.8±1.2	78.8±0.3	83.3±0.1	77.2±0.3	59.4±1.1
CSI (ours)	ResNet-18	99.8±0.0	97.5±0.3	97.6±0.3	90.3±0.3	93.3±0.1	89.2±0.1	79.3±0.2

(b) Unlabeled ImageNet-30									
ImageNet-30 →									
Method	Network	CUB-200	Dogs	Pets	Flowers	Food-101	Places-365	Caltech-256	DTD
Rot [25]	ResNet-18	76.5±0.7	77.2±0.5	70.0±0.5	87.2±0.2	72.7±1.5	52.6±1.4	70.9±0.1	89.9±0.5
Rot+Trans [25]	ResNet-18	74.5±0.5	77.8±1.1	70.0±0.8	86.3±0.3	71.6±1.4	53.1±1.7	70.0±0.2	89.4±0.6
GOAD [2]	ResNet-18	71.5±1.4	74.3±1.6	65.5±1.3	82.8±1.4	68.7±0.7	51.0±1.1	67.4±0.8	87.5±0.8
CSI (ours)	ResNet-18	90.5±0.1	97.1±0.1	85.2±0.2	94.7±0.4	89.2±0.3	78.3±0.3	87.1±0.1	96.9±0.1

Setup. We use ResNet-18 [20] architecture for all the experiments. For data augmentations \mathcal{T} , we adopt those used by Chen et al. [5]: namely, we use the combination of Inception crop [64], horizontal flip, color jitter, and grayscale. For shifting transformations \mathcal{S} , we use the random rotation $0^\circ, 90^\circ, 180^\circ, 270^\circ$ unless specified otherwise, as rotation has the highest OOD-ness (see Section 2.2) values for natural images, *e.g.*, CIFAR-10 [33]. However, we remark that the best shifting transformation can be different for other datasets, *e.g.*, Gaussian noise performs better than rotation for texture datasets (see Table 6 in Section 3.2). By default, we train our models from scratch with the training objective in (5) and detect OOD samples with the ensembled version of the score in (9).

We mainly report the area under the receiver operating characteristic curve (AUROC) as a threshold-free evaluation metric for a detection score. In addition, we report the test accuracy and the expected calibration error (ECE) [45, 19] for the experiments on labeled multi-class datasets. Here, ECE estimates whether a classifier can indicate when they are likely to be incorrect for test samples (from in-distribution) by measuring the difference between prediction confidence and accuracy. The formal description of the metrics and detailed experimental setups are in Appendix A.

3.1 Main results

Unlabeled one-class datasets. We start by considering the *one-class* setup: here, for a given multi-class dataset of C classes, we conduct C one-class classification tasks, where each task chooses one of the classes as in-distribution while the remaining classes being out-of-distribution. We run our experiments on three datasets, following the prior work [15, 25, 2]: CIFAR-10 [33], CIFAR-100 labeled into 20 super-classes [33], and ImageNet-30 [25] datasets. We compare CSI with various prior methods including one-class classifier [59, 56], reconstruction-based [58, 52], and self-supervised [15, 25, 2] approaches. Table 1 summarizes the results, showing that CSI significantly outperforms the prior methods in all the tested cases. We provide the full, additional results, *e.g.*, class-wise AUROC on CIFAR-100 (super-class) and ImageNet-30, in Appendix C.

Unlabeled multi-class datasets. In this setup, we assume that in-distribution samples are from a specific multi-class dataset without labels, testing on various external datasets as out-of-distribution. We compare CSI on two in-distribution datasets: CIFAR-10 [33] and ImageNet-30 [25]. We consider the following datasets as out-of-distribution: SVHN [48], resized LSUN and ImageNet [39], CIFAR-100 [33], and linearly-interpolated samples of CIFAR-10 (Interp.) [11] for CIFAR-10 experiments, and CUB-200 [67], Dogs [29], Pets [51], Flowers [49], Food-101 [3], Places-365 [75], Caltech-256 [18], and DTD [8] for ImageNet-30. We compare CSI with various prior methods, including density-based [11, 55, 61] and self-supervised [15, 2] approaches.

Table 3: Test accuracy (%), ECE (%), and AUROC (%) of confidence-calibrated classifiers trained on labeled (a) CIFAR-10 and (b) ImageNet-30. The reported results are averaged over five trials for CIFAR-10 and one trial for ImageNet-30. Subscripts denote standard deviation, and bold denote the best results. CSI-ens denotes the ensembled prediction, *i.e.*, 4 times slower (as we use rotation).

(a) Labeled CIFAR-10									
Train method	Test acc.	ECE	CIFAR10 →						
			SVHN	LSUN	ImageNet	LSUN (FIX)	ImageNet (FIX)	CIFAR100	Interp.
Cross Entropy	93.0±0.2	6.44±0.2	88.6±0.9	90.7±0.5	88.3±0.6	87.5±0.3	87.4±0.3	85.8±0.3	75.4±0.7
SupCLR [30]	93.8±0.1	5.56±0.1	97.3±0.1	92.8±0.5	91.4±1.2	91.6±1.5	90.5±0.5	88.6±0.2	75.7±0.1
CSI (ours)	94.8±0.1	4.40±0.1	96.5±0.2	96.3±0.5	96.2±0.4	92.1±0.5	92.4±0.0	90.5±0.1	78.5±0.2
CSI-ens (ours)	96.1±0.1	3.50±0.1	97.9±0.1	97.7±0.4	97.6±0.3	93.5±0.4	94.0±0.1	92.2±0.1	80.1±0.3

(b) Labeled ImageNet-30										
Train method	Test acc.	ECE	ImageNet-30 →							
			CUB-200	Dogs	Pets	Flowers	Food-101	Places-365	Caltech-256	DTD
Cross Entropy	94.3	5.08	88.0	96.7	95.0	89.7	79.8	90.5	90.6	90.1
SupCLR [30]	96.9	3.12	86.3	95.6	94.2	92.2	81.2	89.7	90.2	92.1
CSI (ours)	97.0	2.61	93.4	97.7	96.9	96.0	87.0	92.5	91.9	93.7
CSI-ens (ours)	97.8	2.19	94.6	98.3	97.4	96.2	88.9	94.0	93.2	97.4

Table 2 shows the results. Overall, CSI significantly outperforms the prior methods in all benchmarks tested. We remark that CSI is particularly effective for detecting hard (*i.e.*, near-distribution) OOD samples, *e.g.*, CIFAR-100 and Interp. in Table 2a. Also, CSI still shows a notable performance in the cases when prior methods often fail, *e.g.*, AUROC of 50% (*i.e.*, random guess) for Places-365 dataset in Table 2b. Finally, we notice that the resized LSUN and ImageNet datasets officially released by Liang et al. [39] might be misleading to evaluate detection performance for hard OODs: we find that those datasets contain some unintended artifacts, due to incorrect resizing procedure. Such an artifact makes those datasets easily-detectable, *e.g.*, via input statistics. In this respect, we produce and test on their fixed versions, coined LSUN (FIX), and ImageNet (FIX). See Appendix I for details.

Labeled multi-class datasets. We also consider the *labeled* version of the above setting: namely, we now assume that every in-distribution sample also contains discriminative label information. We use the same datasets considered in the unlabeled multi-class setup for in- and out-of-distribution datasets. We train our model as proposed in Section 2.4, and compare it with those trained by other methods, the cross-entropy and supervised contrastive learning (SupCLR) [30]. Since our goal is to calibrate the confidence, the maximum softmax probability is used to detect OOD samples (see [22]).

Table 3 shows the results. Overall, CSI consistently improves AUROC and ECE for all benchmarks tested. Interestingly, CSI also improves test accuracy; even our original purpose of CSI is to learn a representation for OOD detection. CSI can further improve the performance by ensembling over the transformations. We also remark that our results on unlabeled datasets (in Table 2) already show comparable performance to the supervised baselines (in Table 3).

3.2 Ablation study

We perform an ablation study on various shifting transformations, training objectives, and detection scores. Throughout this section, we report the mean AUROC values on one-class CIFAR-10.

Shifting transformation. We measure the OOD-ness (see Section 2.2) of transformations, *i.e.*, the AUROC between in-distribution vs. transformed samples under vanilla SimCLR, and the effects of those transformations when used as a shifting transformation. In particular, we consider Cutout [10], Sobel filtering [28], Gaussian noise, Gaussian blur, and rotation [14]. We remark that these transformations are reported to be ineffective in improving the class discriminative power of SimCLR [5]. We also consider the transformation coined “Perm”, which randomly permutes each part of the evenly partitioned image. Intuitively, such transformations commonly *shift* the input distribution, hence forcing them to be *aligned* can be harmful. Figure 1 visualizes the considered transformations.

Table 4 shows AUROC values of the vanilla SimCLR, where the in-distribution samples shifted by the chosen transformation are given as OOD samples. The shifted samples are easily detected: it validates our intuition that the considered transformations *shift* the input distribution. In particular, “Perm” and “Rotate” are the most distinguishable, which implies they shift the distribution the most.

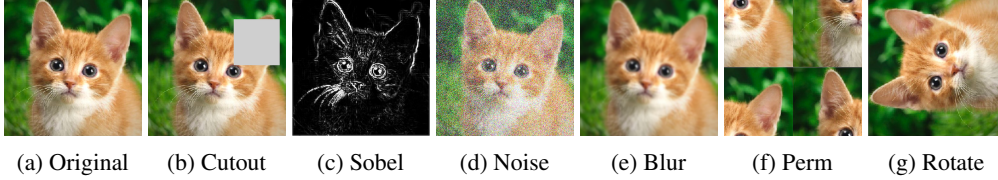


Figure 1: Visualization of the original image and the considered shifting transformations.

Table 4: OOD-ness (%), *i.e.*, the AUROC between in-distribution vs. transformed samples under the vanilla SimCLR (see Section 2.2), of various transformations. The vanilla SimCLR is trained on one-class CIFAR-10 under ResNet-18. Each column denotes the applied transformation.

	Cutout	Sobel	Noise	Blur	Perm	Rotate
OOD-ness	79.5	69.2	74.4	76.0	83.8	85.2

Table 5: Ablation study on various transformations, added or removed from the vanilla SimCLR. “Align” and “Shift” indicates that the transformation is used as \mathcal{T} and \mathcal{S} , respectively. (a) We add a new transformation as an aligned (up) or shifting (down) transformations. (b) We remove (up) or convert-to-shift (down) the transformation from the vanilla SimCLR. All reported values are the mean AUROC (%) over one-class CIFAR-10, and “Base” denotes the vanilla SimCLR.

(a) Add transformations								(b) Remove transformations			
Base		Cutout	Sobel	Noise	Blur	Perm	Rotate		Crop	Jitter	Gray
87.9	+Align	84.3	85.0	85.5	88.0	73.1	76.5	-Align	55.7	78.8	78.4
	+Shift	88.5	88.3	89.3	89.2	90.7	94.3	+Shift	-	-	88.3

Note that “Perm” and “Rotate” turns out to be the most effective shifting transformations; it implies that the transformations *shift* the distribution most indeed performs best for CSI.³

Besides, we apply the transformation upon the vanilla SimCLR: align the transformed samples to the original samples (*i.e.*, use as \mathcal{T}) or consider them as the shifted samples (*i.e.*, use as \mathcal{S}). Table 5a shows that aligning the transformations degrade (or on par) the detection performance, while shifting the transformations gives consistent improvements. We also remove or convert-to-shift the transformation from the vanilla SimCLR in Table 5b, and see similar results. We remark that one can further improve the performance by combining multiple shifting transformations (see Appendix G).

Data-dependence of shifting transformations. We remark that the best shifting transformation depends on the dataset. For example, consider the rotation-invariant datasets: Describable Textures Dataset (DTD) [8] and Textile [60] are in- vs. out-of-distribution, respectively (see Appendix J for more visual examples). For such datasets, rotation (Rot.) does not shift the distribution, and Gaussian noise (Noise) is more suitable transformation (see Table 6a). Table 6b shows that CSI using Gaussian noise (“CSI(N)”) indeed improves the vanilla SimCLR (“Base”) while CSI using rotation (“CSI(R)”) degrades instead. This results support our principles on selecting shifting transformations.

Table 6: OOD-ness (%) and AUROC (%) on DTD, where Textile is used for OOD.

(a) OOD-ness				(b) AUROC		
Rot.	Noise			Base	CSI(R)	CSI(N)
50.6	75.7			70.3	65.9	80.1

Linear evaluation. We also measure the linear evaluation [32], the accuracy of a linear classifier to discriminate classes of in-distribution samples. It is widely used for evaluating the quality of (unsupervised) learned representation. We report the linear evaluation of vanilla SimCLR and CSI (with shifting rotation), trained under unlabeled CIFAR-10. They show comparable results, 90.48% for SimCLR and 90.19% for CSI; CSI is more specialized to learn a representation for OOD detection.

Training objective. In Table 7a, we assess the individual effects of each component that consists of our final training objective (5): namely, we compare the vanilla SimCLR (2), contrasting shifted

³We also have tried contrasting *external* OOD samples similarly to [24]; however, we find that naïvely using them in our framework degrade the performance. This is because the contrastive loss also discriminates *within* external OOD samples, which is unnecessary and an additional learning burden for our purpose.

Table 7: Ablation study on each component of our proposed (a) training objective and (b) detection score. For (a), we use the corresponding detection score for each training loss; namely, (6) to (9) for (2) to (5), respectively. For (b), we use the model trained by the final training loss (5). We measure the mean AUROC (%) values, trained under CIFAR-10 with ResNet-18. Each row indicates the corresponding equation of the given checkmarks, and bold denotes the best results. “Con.”, “Cls.”, and “Ensem.” denotes contrast, classify, and ensemble, respectively.

	(a) Training objective					(b) Detection score			
	SimCLR	Con.	Cls.	AUROC		Con.	Cls.	Ensem.	AUROC
$\mathcal{L}_{\text{SimCLR}}$ (2)	✓	-	-	87.9	s_{con} (6)	✓	-	-	91.3
$\mathcal{L}_{\text{con-SI}}$ (3)	✓	✓	-	91.6	$s_{\text{con-SI}}$ (7)	✓	-	✓	93.3
$\mathcal{L}_{\text{cls-SI}}$ (4)	-	-	✓	88.6	$s_{\text{cls-SI}}$ (8)	-	✓	✓	93.8
\mathcal{L}_{CSI} (5)	✓	✓	✓	94.3	s_{CSI} (9)	✓	✓	✓	94.3

instances (3), and classifying shifted instances (4) losses. For the evaluation of the models of different training objectives (2) to (5), we use the detection scores defined in (6) to (9), respectively. We remark that both contrasting and classifying shows better results than the vanilla SimCLR; and combining them (*i.e.*, the final CSI objective (5)) gives further improvements, *i.e.*, two losses are complementary.

Detection score. Finally, Table 7b shows the effect of each component in our detection score: the vanilla contrastive (6), contrasting shifted instances (7), and classifying shifted instances (8) scores. We ensemble the scores over both \mathcal{T} and \mathcal{S} for (7) to (9), and use a single sample for (6). All the reported values are evaluated from the model trained by the final objective 5. Similar to above, both contrasting and classifying scores show better results than the vanilla contrastive score; and combining them (*i.e.*, the final CSI score (9)) gives further improvements.

4 Related work

OOD detection. Recent works on unsupervised OOD detection (*i.e.*, no external OOD samples) [26] can be categorized as: (a) density-based [74, 46, 6, 47, 11, 55, 61, 17], (b) reconstruction-based [58, 76, 9, 54, 52, 7], (c) one-class classifier [59, 56], and (d) self-supervised [15, 25, 2] approaches. Our work falls into (c) the self-supervised approach, as it utilizes the representation learned from self-supervision [14]. However, unlike prior works [15, 25, 2] focusing on the self-label classification tasks (*e.g.*, predict the angle of the rotated image), we first incorporate *contrastive learning* [5] for OOD detection. Concurrently, Winkens et al. [68] and Liu and Abbeel [40] report that contrastive learning also improves the OOD detection performance of classifiers [39, 38, 25].

Confidence-calibrated classifiers. Confidence-calibrated classifiers aim to calibrate the prediction confidence (maximum softmax probability), which can be directly used as an uncertainty estimator for both within in-distribution [45, 19] and in- vs. out-of-distribution [22, 37]. Prior works improved calibration through inference [19] or training [37] schemes, which are can be jointed applied to our method. Some works design a specific detection score upon the pre-trained classifiers [39, 38], but they only target OOD detection, while ours also consider the in-distribution calibration.

Self-supervised learning. Self-supervised learning [14, 32], particularly contrastive learning [13] via instance discrimination [69], has shown remarkable success on visual representation learning [21, 5]. However, most prior works focus on the downstream tasks (*e.g.*, classification), and other advantages (*e.g.*, uncertainty or robustness) are rarely investigated [25, 31]. Our work, concurrent with [40, 68], first verifies that contrastive learning is also effective for OOD detection. In particular, we find that the shifting transformations, which were known to be harmful and unused for the standard contrastive learning [5], can help OOD detection. This observation provides new considerations for selecting transformations, *i.e.*, which transformation should be used for positive or negative [66, 71].

We further provide a more comprehensive survey and discussions with prior works in Appendix B.

5 Conclusion

We propose a simple yet effective method named contrasting shifted instances (CSI), which extends the power of contrastive learning for out-of-distribution (OOD) detection problems. CSI demonstrates outstanding performance under various OOD detection scenarios. We believe our work would guide various future directions in OOD detection and self-supervised learning as an important baseline.

Acknowledgements

This work was supported by Institute of Information & Communications Technology Planning & Evaluation (IITP) grant funded by the Korea government (MSIT) (No.2019-0-00075, Artificial Intelligence Graduate School Program (KAIST) and No.2017-0-01779, A machine learning and statistical inference framework for explainable artificial intelligence). We thank Sihyun Yu, Chaewon Kim, Hyuntak Cha, Hyunwoo Kang, and Seunghyun Lee for helpful feedback and suggestions.

Broader Impact

This paper is focused on the subject of *out-of-distribution (OOD)* (or novelty, anomaly) detection, which is an essential ingredient for building safe and reliable intelligent systems [1]. We expect our results to have two consequences for academia and broader society.

Rethinking representation for OOD detection. In this paper, we demonstrate that the representation for classification (or other related tasks, measured by linear evaluation [32]) can be different from the representation for OOD detection. In particular, we verify that the “hard” augmentations, thought to be harmful for contrastive representation learning [5], can be helpful for OOD detection. Our observation raises new questions for both representation learning and OOD detection: (a) representation learning researches should also report the OOD detection results as an evaluation metric, (b) OOD detection researches should more investigate the specialized representation.

Towards reliable intelligent system. The intelligent system should be robust to the potential dangers of uncertain environments (*e.g.*, financial crisis [65]) or malicious adversaries (*e.g.*, cybersecurity [34]). Detecting outliers is also related to human safety (*e.g.*, medical diagnosis [4] or autonomous driving [12]), and has a broad range of industrial applications (*e.g.*, manufacturing inspection [42]). However, the system can be stuck into *confirmation bias*, *i.e.*, ignore new information with a myopic perspective. We hope the system to balance the exploration and exploitation of the knowledge.

References

- [1] D. Amodei, C. Olah, J. Steinhardt, P. Christiano, J. Schulman, and D. Mané. Concrete problems in ai safety. *arXiv preprint arXiv:1606.06565*, 2016.
- [2] L. Bergman and Y. Hoshen. Classification-based anomaly detection for general data. In *International Conference on Learning Representations*, 2020.
- [3] L. Bossard, M. Guillaumin, and L. Van Gool. Food-101—mining discriminative components with random forests. In *European Conference on Computer Vision*, 2014.
- [4] R. Caruana, Y. Lou, J. Gehrke, P. Koch, M. Sturm, and N. Elhadad. Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day readmission. In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2015.
- [5] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton. A simple framework for contrastive learning of visual representations. In *International Conference on Machine Learning*, 2020.
- [6] H. Choi, E. Jang, and A. A. Alemi. Waic, but why? generative ensembles for robust anomaly detection. *arXiv preprint arXiv:1810.01392*, 2018.
- [7] S. Choi and S.-Y. Chung. Novelty detection via blurring. In *International Conference on Learning Representations*, 2020.
- [8] M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, and A. Vedaldi. Describing textures in the wild. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2014.
- [9] L. Deecke, R. Vandermeulen, L. Ruff, S. Mandt, and M. Kloft. Image anomaly detection with generative adversarial networks. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, 2018.
- [10] T. DeVries and G. W. Taylor. Improved regularization of convolutional neural networks with cutout. *arXiv preprint arXiv:1708.04552*, 2017.

- [11] Y. Du and I. Mordatch. Implicit generation and modeling with energy based models. In *Advances in Neural Information Processing Systems*, 2019.
- [12] K. Eykholt, I. Evtimov, E. Fernandes, B. Li, A. Rahmati, C. Xiao, A. Prakash, T. Kohno, and D. Song. Robust physical-world attacks on deep learning visual classification. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
- [13] W. Falcon and K. Cho. A framework for contrastive self-supervised learning and designing a new approach. *arXiv preprint arXiv:2009.00104*, 2020.
- [14] S. Gidaris, P. Singh, and N. Komodakis. Unsupervised representation learning by predicting image rotations. In *International Conference on Learning Representations*, 2018.
- [15] I. Golan and R. El-Yaniv. Deep anomaly detection using geometric transformations. In *Advances in Neural Information Processing Systems*, 2018.
- [16] P. Goyal, P. Dollár, R. Girshick, P. Noordhuis, L. Wesolowski, A. Kyrola, A. Tulloch, Y. Jia, and K. He. Accurate, large minibatch sgd: Training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017.
- [17] W. Grathwohl, K.-C. Wang, J.-H. Jacobsen, D. Duvenaud, M. Norouzi, and K. Swersky. Your classifier is secretly an energy based model and you should treat it like one. In *International Conference on Learning Representations*, 2020.
- [18] G. Griffin, A. Holub, and P. Perona. Caltech-256 object category dataset, 2007.
- [19] C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger. On calibration of modern neural networks. In *International Conference on Machine Learning*, 2017.
- [20] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- [21] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick. Momentum contrast for unsupervised visual representation learning. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2020.
- [22] D. Hendrycks and K. Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. In *International Conference on Learning Representations*, 2017.
- [23] D. Hendrycks, K. Lee, and M. Mazeika. Using pre-training can improve model robustness and uncertainty. In *International Conference on Machine Learning*, 2019.
- [24] D. Hendrycks, M. Mazeika, and T. Dietterich. Deep anomaly detection with outlier exposure. In *International Conference on Learning Representations*, 2019.
- [25] D. Hendrycks, M. Mazeika, S. Kadavath, and D. Song. Using self-supervised learning can improve model robustness and uncertainty. In *Advances in Neural Information Processing Systems*, 2019.
- [26] V. Hodge and J. Austin. A survey of outlier detection methodologies. *Artificial intelligence review*, 2004.
- [27] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *International Conference on Machine Learning*, 2015.
- [28] N. Kanopoulos, N. Vasanthavada, and R. L. Baker. Design of an image edge detection filter using the sobel operator. *IEEE Journal of solid-state circuits*, 1988.
- [29] A. Khosla, N. Jayadevaprakash, B. Yao, and L. Fei-Fei. Novel dataset for fine-grained image categorization. In *IEEE Conference on Computer Vision and Pattern Recognition*, Colorado Springs, CO, June 2011.
- [30] P. Khosla, P. Teterwak, C. Wang, A. Sarna, Y. Tian, P. Isola, A. Maschinot, C. Liu, and D. Krishnan. Supervised contrastive learning. In *Advances in Neural Information Processing Systems*, 2020.

- [31] M. Kim, J. Tack, and S. J. Hwang. Adversarial self-supervised contrastive learning. In *Advances in Neural Information Processing Systems*, 2020.
- [32] A. Kolesnikov, X. Zhai, and L. Beyer. Revisiting self-supervised visual representation learning. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [33] A. Krizhevsky et al. Learning multiple layers of features from tiny images, 2009.
- [34] C. Kruegel and G. Vigna. Anomaly detection of web-based attacks. In *Proceedings of the 10th ACM conference on Computer and communications security*, 2003.
- [35] H. Lee, S. J. Hwang, and J. Shin. Self-supervised label augmentation via input transformations. In *International Conference on Machine Learning*, 2020.
- [36] K. Lee, C. Hwang, K. S. Park, and J. Shin. Confident multiple choice learning. In *International Conference on Machine Learning*, 2017.
- [37] K. Lee, H. Lee, K. Lee, and J. Shin. Training confidence-calibrated classifiers for detecting out-of-distribution samples. In *International Conference on Learning Representations*, 2018.
- [38] K. Lee, K. Lee, H. Lee, and J. Shin. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. In *Advances in Neural Information Processing Systems*, 2018.
- [39] S. Liang, Y. Li, and R. Srikant. Enhancing the reliability of out-of-distribution image detection in neural networks. In *International Conference on Learning Representations*, 2018.
- [40] H. Liu and P. Abbeel. Hybrid discriminative-generative training via contrastive learning. *arXiv preprint arXiv:2007.09070*, 2020.
- [41] I. Loshchilov and F. Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016.
- [42] D. Lucke, C. Constantinescu, and E. Westkämper. Smart factory-a step towards the next generation of manufacturing. *Manufacturing Systems and Technologies for the New Frontier*, page 115, 2008.
- [43] L. v. d. Maaten and G. Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 2008.
- [44] J. MacQueen et al. Some methods for classification and analysis of multivariate observations, 1967.
- [45] M. P. Naeini, G. Cooper, and M. Hauskrecht. Obtaining well calibrated probabilities using bayesian binning. In *AAAI Conference on Artificial Intelligence*, 2015.
- [46] E. Nalisnick, A. Matsukawa, Y. W. Teh, D. Gorur, and B. Lakshminarayanan. Do deep generative models know what they don’t know? In *International Conference on Learning Representations*, 2019.
- [47] E. Nalisnick, A. Matsukawa, Y. W. Teh, and B. Lakshminarayanan. Detecting out-of-distribution inputs to deep generative models using a test for typicality. *arXiv preprint arXiv:1906.02994*, 2019.
- [48] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng. Reading digits in natural images with unsupervised feature learning. In *Advances in Neural Information Processing Systems Workshop on Deep Learning and Unsupervised Feature Learning*, 2011.
- [49] M.-E. Nilsback and A. Zisserman. A visual vocabulary for flower classification. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2006.
- [50] A. v. d. Oord, Y. Li, and O. Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- [51] O. M. Parkhi, A. Vedaldi, A. Zisserman, and C. Jawahar. Cats and dogs. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2012.

- [52] P. Perera, R. Nallapati, and B. Xiang. Ogan: One-class novelty detection using gans with constrained latent representations. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [53] C. Phua, V. Lee, K. Smith, and R. Gayler. A comprehensive survey of data mining-based fraud detection research. *arXiv preprint arXiv:1009.6119*, 2010.
- [54] S. Pidhorskyi, R. Almhosen, and G. Doretto. Generative probabilistic novelty detection with adversarial autoencoders. In *Advances in Neural Information Processing Systems*, 2018.
- [55] J. Ren, P. J. Liu, E. Fertig, J. Snoek, R. Poplin, M. Depristo, J. Dillon, and B. Lakshminarayanan. Likelihood ratios for out-of-distribution detection. In *Advances in Neural Information Processing Systems*, 2019.
- [56] L. Ruff, R. Vandermeulen, N. Goernitz, L. Deecke, S. A. Siddiqui, A. Binder, E. Müller, and M. Kloft. Deep one-class classification. In *International Conference on Machine Learning*, 2018.
- [57] L. Ruff, R. A. Vandermeulen, N. Görnitz, A. Binder, E. Müller, K.-R. Müller, and M. Kloft. Deep semi-supervised anomaly detection. In *International Conference on Learning Representations*, 2020.
- [58] T. Schlegl, P. Seeböck, S. M. Waldstein, U. Schmidt-Erfurth, and G. Langs. Unsupervised anomaly detection with generative adversarial networks to guide marker discovery. In *International conference on information processing in medical imaging*, 2017.
- [59] B. Schölkopf, R. C. Williamson, A. J. Smola, J. Shawe-Taylor, and J. C. Platt. Support vector method for novelty detection. In *Advances in Neural Information Processing Systems*, 2000.
- [60] H. Schulz-Mirbach. Tilda-ein referenzdatensatz zur evaluierung von sichtprüfungsverfahren für textiloberflächen. *Interner Bericht*, 1996. URL <https://lmb.informatik.uni-freiburg.de/resources/datasets/tilda.en.html>.
- [61] J. Serrà, D. Álvarez, V. Gómez, O. Slizovskaia, J. F. Núñez, and J. Luque. Input complexity and out-of-distribution detection with likelihood-based generative models. In *International Conference on Learning Representations*, 2020.
- [62] Severstal. Severstal: Steel defect detection, 2019. URL <https://www.kaggle.com/c/severstal-steel-defect-detection>.
- [63] A. Srinivas, M. Laskin, and P. Abbeel. Curl: Contrastive unsupervised representations for reinforcement learning. In *International Conference on Machine Learning*, 2020.
- [64] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–9, 2015.
- [65] J. B. Taylor and J. C. Williams. A black swan in the money market. *American Economic Journal: Macroeconomics*, 2009.
- [66] Y. Tian, C. Sun, B. Poole, D. Krishnan, C. Schmid, and P. Isola. What makes for good views for contrastive learning. In *Advances in Neural Information Processing Systems*, 2020.
- [67] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The Caltech-UCSD Birds-200-2011 Dataset. Technical Report CNS-TR-2011-001, California Institute of Technology, 2011.
- [68] J. Winkens, R. Bunel, A. G. Roy, R. Stanforth, V. Natarajan, J. R. Ledsam, P. MacWilliams, P. Kohli, A. Karthikesalingam, and S. Kohl. Contrastive training for improved out-of-distribution detection. *arXiv preprint arXiv:2007.05566*, 2020.
- [69] Z. Wu, Y. Xiong, S. X. Yu, and D. Lin. Unsupervised feature learning via non-parametric instance discrimination. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018.

- [70] G.-S. Xia, X. Bai, J. Ding, Z. Zhu, S. Belongie, J. Luo, M. Datcu, M. Pelillo, and L. Zhang. Dota: A large-scale dataset for object detection in aerial images. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
- [71] T. Xiao, X. Wang, A. A. Efros, and T. Darrell. What should not be contrastive in contrastive learning. *arXiv preprint arXiv:2008.05659*, 2020.
- [72] Y. You, I. Gitman, and B. Ginsburg. Large batch training of convolutional networks. *arXiv preprint arXiv:1708.03888*, 2017.
- [73] S. Yun, J. Park, K. Lee, and J. Shin. Regularizing class-wise predictions via self-knowledge distillation. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2020.
- [74] S. Zhai, Y. Cheng, W. Lu, and Z. Zhang. Deep structured energy based models for anomaly detection. In *International Conference on Machine Learning*, 2016.
- [75] B. Zhou, A. Lapedriza, A. Khosla, A. Oliva, and A. Torralba. Places: A 10 million image database for scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017.
- [76] B. Zong, Q. Song, M. R. Min, W. Cheng, C. Lumezanu, D. Cho, and H. Chen. Deep autoencoding gaussian mixture model for unsupervised anomaly detection. In *International Conference on Learning Representations*, 2018.

Appendix

CSI: Novelty Detection via Contrastive Learning on Distributionally Shifted Instances

A Experimental details

Training details. We use ResNet-18 [20] as the base encoder network f_θ and 2-layer multi-layer perceptron with 128 embedding dimension as the projection head g_ϕ . All models are trained by minimizing the final loss \mathcal{L}_{CSI} (5) with a temperature of $\tau = 0.5$. We follow the same optimization step of SimCLR [5]. For optimization, we train CSI with 1,000 epoch under LARS optimizer [72] with weight decay of $1\text{e-}6$ and momentum with 0.9. For the learning rate scheduling, we use linear warmup [16] for early 10 epochs until learning rate of 1.0 and decay with cosine decay schedule without a restart [41]. We use batch size of 512 for both vanilla SimCLR and ours: where the batch is given by \mathcal{B} for vanilla SimCLR and the aggregated one $\bigcup_{S \in \mathcal{S}} \mathcal{B}_S$ for ours. Furthermore, we use global batch normalization (BN) [27], which shares the BN parameters (mean and variance) over the GPUs in distributed training.

For supervised contrastive learning (SupCLR) [30] and supervised CSI, we select the best temperature from $\{0.07, 0.5\}$: SupCLR recommend 0.07 but 0.5 was better in our experiments. For training the encoder f_θ , we use the same optimization scheme as above, except using 700 for the epoch. For training the linear classifier, we train the model for 100 epochs with batch size 128, using stochastic gradient descent with momentum 0.9. The learning rate starts at 0.1 and is dropped by a factor of 10 at 60%, 75%, and 90% of the training progress.

Data augmentation details. We use SimCLR augmentations: Inception crop [64], horizontal flip, color jitter, and grayscale for random augmentations \mathcal{T} , and rotation as shifting transformation \mathcal{S} . The detailed description of the augmentations are as follows:

- **Inception crop.** Randomly crops the area of the original image with uniform distribution 0.08 to 1.0. After the crop, cropped image are resized to the original image size.
- **Horizontal flip.** Flips the image horizontally with 50% of probability.
- **Color jitter.** Change the hue, brightness, and saturation of the image. We transform the RGB (red, green, blue) image into an HSV (hue, saturation, value) image format and add noise to the HSV channels. We apply color jitter with 80% of probability.
- **Grayscale.** Convert into a gray image. Randomly apply a grayscale with 20% of probability.
- **Rotation.** We use rotation as \mathcal{S} , the shifting transformation, $\{0^\circ, 90^\circ, 180^\circ, 270^\circ\}$. For a given batch \mathcal{B} , we apply each rotation degree to obtain the new batch for CSI: $\bigcup_{S \in \mathcal{S}} \mathcal{B}_S$.

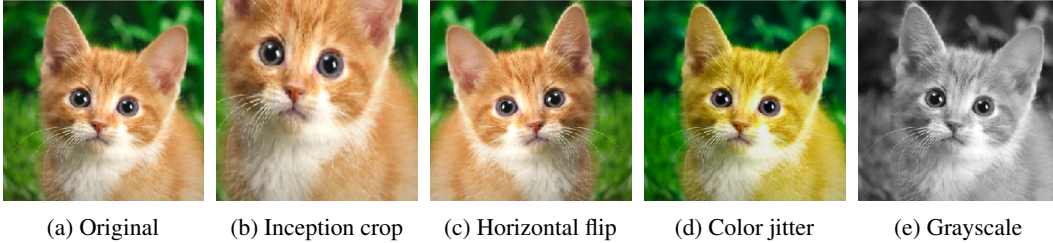


Figure 2: Visualization of original image and SimCLR augmentations.

Dataset details. For one-class datasets, we train one class of CIFAR-10 [33], CIFAR-100 (super-class) [33], and ImageNet-30 [25]. CIFAR-10 and CIFAR-100 consist of 50,000 training and 10,000 test images with 10 and 20 (super-class) image classes, respectively. ImageNet-30 contains 39,000 training and 3,000 test images with 30 image classes.

For unlabeled and labeled multi-class datasets, we train ResNet with CIFAR-10 and ImageNet-30. For CIFAR-10, out-of-distribution (OOD) samples are as follows: SVHN [48] consists of 26,032 test images with 10 digits, resized LSUN [39] consists of 10,000 test images of 10 different scenes, resized ImageNet [39] consists of 10,000 test images with 200 images classes from a subset of full ImageNet dataset, Interp. consists of 10,000 test images of linear interpolation of CIFAR-10 test images, and LSUN (FIX), ImageNet (FIX) consists of 10,000 test images, respectively with following details in Appendix I. For multi-class ImageNet-30, OOD samples are as follows: CUB-200 [67], Stanford Dogs [29], Oxford Pets [51], Oxford Flowers [49], Food-101 [3] without the “hotdog” class to avoid overlap, Places-365 [75] with small images (256 * 256) validation set, Caltech-256 [18], and Describable Textures Dataset (DTD) [8]. Here, we randomly sample 3,000 images to balance with the in-distribution test set.

Evaluation metrics. For evaluation, we measure the two metrics that each measures (a) the effectiveness of the proposed score in distinguishing in- and out-of-distribution images, (b) the confidence calibration of softmax classifier.

- **Area under the receiver operating characteristic curve (AUROC).** Let TP, TN, FP, and FN denote true positive, true negative, false positive and false negative, respectively. The ROC curve is a graph plotting true positive rate = $TP / (TP+FN)$ against the false positive rate = $FP / (FP+TN)$ by varying a threshold.
- **Expected calibration error (ECE).** For a given test data $\{(x_n, y_n)\}_{n=1}^N$, we group the predictions into M interval bins (each of size $1/M$). Let B_m be the set of indices of samples whose prediction confidence falls into the interval $(\frac{m-1}{M}, \frac{m}{M}]$. Then, the expected calibration error (ECE) [45, 19] is follows:

$$ECE = \sum_{m=1}^M \frac{|B_m|}{N} |\text{acc}(B_m) - \text{conf}(B_m)|, \quad (13)$$

where $\text{acc}(B_m)$ is accuracy of B_m : $\text{acc}(B_m) = \frac{1}{|B_m|} \sum_{i \in B_m} \mathbb{1}_{\{y_i = \arg \max_y p(y|x_i)\}}$ where $\mathbb{1}$ is indicator function and $\text{conf}(B_m)$ is confidence of B_m : $\text{conf}(B_m) = \frac{1}{|B_m|} \sum_{i \in B_m} q(x_i)$ where $q(x_i)$ is the confidence of data x_i .

B Detailed review on related work

B.1 OOD detection

Out-of-distribution (OOD) detection is a classic and essential problem in machine learning, studied under different names, *e.g.*, novelty or anomaly detection [26]. In this paper, we primarily focus on *unsupervised* OOD detection, which is arguably the most traditional and popular setup in the field [59]. In this setting, the detector can only access in-distribution samples while required to identify unseen OOD samples. There are other settings, *e.g.*, semi-supervised setting - the detector can access a small subset of out-of-distribution samples [24, 57], or supervised setting - the detector knows the target out-of-distribution, but we do not consider those settings in this paper. We remark that the unsupervised setting is the most practical and challenging scenario since there are *infinitely* many cases for out-of-distribution, and it is often not possible to have such external data.

Most recent works can be categorized as: (a) density-based [74, 46, 6, 47, 11, 55, 61, 17], (b) reconstruction-based [58, 76, 9, 54, 52, 7], (c) one-class classifier [59, 56, 57], and (d) self-supervised [15, 25, 2] approaches. We note that there are more extensive literature on this topic, but we mainly focus on the recent work based on deep learning. Brief description for each method are as follows:

- **Density-based methods.** Density-based methods are one of the most classic and principled approaches for OOD detection. Intuitively, they directly use the likelihood of the sample as the detection score. However, recent studies reveal that the likelihood is often not the best metric - especially for deep neural networks with complex datasets [46]. Several work thus proposed modified scores, *e.g.*, typicality [47], WAIC [6], likelihood ratio [55], input complexity [61], or unnormalized likelihood (*i.e.*, energy) [11, 17].
- **Reconstruction-based methods.** Reconstruction-based approach is another popular line of research for OOD detection. It trains an encoder-decoder network that reconstructs the training data in an unsupervised manner. Since the network would less generalize for unseen OOD samples, they use the reconstruction loss as a detection score. Some works utilize auto-encoders [76, 54] or generative adversarial networks [58, 9, 52].
- **One-class classifiers.** One-class classifiers are also a classic and principled approach for OOD detection. They learn a decision boundary of in- vs. out-of-distribution samples by giving some margin covering the in-distribution samples [59]. Recent works have shown that the one-class classifier is effective upon the deep representation [56].
- **Self-supervised methods.** Self-supervised approaches are a relatively new technique based on the rich representation learned from self-supervision [14]. They train a network with a pre-defined task (*e.g.*, predict the angle of the rotated image) on the training set, and use the generalization error to detect OOD samples. Recent self-supervised approaches show outstanding results on various OOD detection benchmark datasets [15, 25, 2].

Our work falls into (c) the self-supervised approach [15, 25, 2]. However, unlike prior work focusing on the self-label classification tasks (*e.g.*, rotation [14]) which trains an auxiliary classifier to predict the transformation applied to the sample, we first incorporate *contrastive learning* [5] for OOD detection. To that end, we design a novel detection score utilizing the unique characteristic of contrastive learning, *e.g.*, the features in the projection layer learned by cosine similarity. We also propose a novel self-supervised training scheme that further improves the representation for OOD detection. Nevertheless, we acknowledge that the prior work largely inspired our work. For instance, the classifying shifted instances loss (4) follows the form of auxiliary classifiers [25], which gives further improvement upon our novel contrasting shifted instances loss (3).

Concurrently, Winkens et al. [68] and Liu and Abbeel [40] report the similar observations that contrastive learning also improves the OOD detection performance of classifiers [39, 38, 25]. Winkens et al. [68] jointly train a classifier with the SimCLR [5] objective and use the Mahalanobis distance [38] as a detection score. Liu and Abbeel [40] approximates JEM [17] (a joint model of classifier and energy-based model [11]) by a combination of classification and contrastive loss and use density-based detection scores [17]. In contrast to both work, we mainly focus on the *unlabeled OOD* setting (although we also discuss the confident-calibrated classifiers). Here, we design a novel detection score, since how to utilize the contrastive representation (which is learned in an unsupervised manner) for OOD detection have not been explored before.

B.2 Confidence-calibrated classifiers

Another line of research is on confidence-calibrated classifiers [22], which relaxes the overconfidence issue of the classifiers. There are two types of calibration: (a) *in-distribution* calibration [45, 19], that aligns the uncertainty and the actual accuracy, measured by ECE, and (b) *out-of-distribution* detection [22, 37], that reduces the uncertainty of OOD samples, measured by AUROC. Note that the goal of confidence-calibrated classifiers is to regularize the prediction. Hence, the softmax probability is used for all three tasks: classification, in-distribution calibration, and out-of-distribution detection. Namely, the detection score is given by the prediction confidence (or maximum softmax probability) [22]. Prior works improved calibration through inference (temperature scaling) [19] or training (regularize predictions of OOD samples) [37] schemes, which can be jointly applied to our method. Some works design a specific detection score upon the pre-trained classifiers [39, 38], but they only target OOD detection, while ours also consider the in-distribution calibration.

B.3 Self-supervised learning

Self-supervised learning [14, 32] has shown remarkable success in learning representations. In particular, contrastive learning [13] via instance discrimination [69] show the state-of-the-art results on visual representation learning [21, 5]. However, most prior works focus on improving the downstream task performance (*e.g.*, classification), and other advantages of self-supervised learning (*e.g.*, uncertainty or robustness) are rarely investigated [25, 31]. Our work, concurrent with [40, 68], first verifies that contrastive learning is also effective for OOD detection.

Furthermore, we find that the shifting transformations, which were known to be harmful and unused for the standard contrastive learning [5], can help OOD detection. This observation provides new considerations for selecting transformations, *i.e.*, which transformation should be used for positive or negative [66, 71]. Specifically, Tian et al. [66] claims the optimal views (or transformations) of the *positive* pairs should minimize the mutual information while keeping the task-relevant information. It suggests that the shifting transformation may not contain the information for classification, but may contain OOD detection information when used for the *negative* pairs. Xiao et al. [71] suggests a framework that automatically learns whether the transformation should be positive or negative. One could consider incorporating our principle on shifting transformation (*i.e.*, OOD-ness); OOD detection could be another evaluation metric for the learned representations.

C Additional one-class OOD detection results

Table 8 presents the confusion matrix of AUROC values of our method on one-class CIFAR-10 datasets, where bold denotes the hard pairs. The results align with the human intuition that ‘car’ is confused to ‘ship’ and ‘truck’, and ‘cat’ is confused to ‘dog’.

Table 9 presents the OOD detection results of various methods on one-class CIFAR-100 (super-class) datasets, for all 20 super-classes. Our method outperforms the prior methods for all classes.

Table 10 presents the OOD detection results of our method on one-class ImageNet-30 dataset, for all 30 classes. Our method consistently performs well for all classes.

Table 8: Confusion matrix of AUROC (%) values of our method on one-class CIFAR-10. The row and column indicates the in-distribution and OOD class, respectively, and the final column indicates the mean value. Bold denotes the values under 80%, which implies the hard pair.

	Plane	Car	Bird	Cat	Deer	Dog	Frog	Horse	Ship	Truck	Mean
Plane	-	74.1	95.8	98.4	94.9	98.0	96.2	90.1	79.6	82.8	90.0
Car	99.3	-	99.9	99.9	99.8	99.9	99.8	99.7	98.7	95.0	99.1
Bird	91.1	97.5	-	97.3	87.0	92.5	96.1	83.2	96.4	98.0	93.2
Cat	91.9	91.5	90.3	-	83.3	67.0	89.6	79.0	92.8	91.9	86.4
Deer	95.7	98.4	94.9	96.6	-	94.7	98.7	69.0	97.4	98.8	93.8
Dog	97.9	98.5	95.5	90.3	88.1	-	96.8	76.6	98.6	98.3	93.4
Frog	93.6	92.3	94.6	96.1	96.8	96.3	-	95.2	94.4	97.3	95.2
Horse	99.3	99.5	99.0	99.3	94.2	97.4	99.8	-	99.7	99.4	98.6
Ship	96.6	91.2	99.5	99.7	99.4	99.7	99.5	99.3	-	96.6	97.9
Truck	96.2	72.3	99.4	99.5	99.1	99.4	98.7	98.3	96.2	-	95.5

Table 9: AUROC (%) values of various OOD detection methods trained on one-class CIFAR-100 (super-class). Each row indicates the results of the selected super-class, and the final row indicates the mean value. * denotes the values from the reference, and bold denotes the best results.

	OC-SVM*	DAGMM*	DSEBM*	ADGAN*	Geom*	Rot	Rot+Trans	GOAD	CSI (ours)
0	68.4	43.4	64.0	63.1	74.7	78.6	79.6	73.9	86.3
1	63.6	49.5	47.9	64.9	68.5	73.4	73.3	69.2	84.8
2	52.0	66.1	53.7	41.3	74.0	70.1	71.3	67.6	88.9
3	64.7	52.6	48.4	50.0	81.0	68.6	73.9	71.8	85.7
4	58.2	56.9	59.7	40.6	78.4	78.7	79.7	72.7	93.7
5	54.9	52.4	46.6	42.8	59.1	69.7	72.6	67.0	81.9
6	57.2	55.0	51.7	51.1	81.8	78.8	85.1	80.0	91.8
7	62.9	52.8	54.8	55.4	65.0	62.5	66.8	59.1	83.9
8	65.6	53.2	66.7	59.2	85.5	84.2	86.0	79.5	91.6
9	74.1	42.5	71.2	62.7	90.6	86.3	87.3	83.7	95.0
10	84.1	52.7	78.3	79.8	87.6	87.1	88.6	84.0	94.0
11	58.0	46.4	62.7	53.7	83.9	76.2	77.1	68.7	90.1
12	68.5	42.7	66.8	58.9	83.2	83.3	84.6	75.1	90.3
13	64.6	45.4	52.6	57.4	58.0	60.7	62.1	56.6	81.5
14	51.2	57.2	44.0	39.4	92.1	87.1	88.0	83.8	94.4
15	62.8	48.8	56.8	55.6	68.3	69.0	71.9	66.9	85.6
16	66.6	54.4	63.1	63.3	73.5	71.7	75.6	67.5	83.0
17	73.7	36.4	73.0	66.7	93.8	92.2	93.5	91.6	97.5
18	52.8	52.4	57.7	44.3	90.7	90.4	91.5	88.0	95.9
19	58.4	50.3	55.5	53.0	85.0	86.5	88.1	82.6	95.2
Mean	63.1	50.6	58.8	55.2	78.7	77.7	79.8	74.5	89.6

Table 10: AUROC (%) values of our method trained on one-class ImageNet-30. The first and third row indicates the selected class, and the second and fifth row indicates the corresponding results.

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
85.9	99.0	99.8	90.5	95.8	99.2	96.6	83.5	92.2	84.3	99.0	94.5	97.1	87.7	96.4
15	16	17	18	19	20	21	22	23	24	25	26	27	28	29
84.7	99.7	75.6	95.2	73.8	94.7	95.2	99.2	98.5	82.5	89.7	82.1	97.2	82.1	97.6

D Ablation study on random augmentation

We verify that ensembling the scores over the random augmentations \mathcal{T} improves OOD detection. However, naïve random sampling from the entire \mathcal{T} is often sample inefficient. We find that choosing a proper subset $\mathcal{T}_{\text{control}} \subset \mathcal{T}$ improves the performance for given number of samples. Specifically, we choose $\mathcal{T}_{\text{control}}$ as the set of the *most common* samples. For example, the size of the cropping area is sampled from $\mathcal{U}[0.08, 1]$ for uniform distribution \mathcal{U} during training. Since the rare samples, *e.g.*, area near 0.08 increases the noise, we only use the samples with size $(0.08 + 1)/2 = 0.54$ during inference. Table 11 shows random sampling from the controlled set often gives improvements.

Table 11: AUROC (%) values of our method for different number of random augmentations, under one-class (OC-) CIFAR-10 and CIFAR-100 (super-class). The values are averaged over classes. Random augmentations over the controlled set show the best performance.

# of samples	Controlled	OC-CIFAR-10	OC-CIFAR-100
4	-	92.22	87.36
40	-	94.13	89.51
40	✓	94.31	89.55

E Efficient computation of (6) via coresets

One can reduce the computation and memory cost of the contrastive score (6) by selecting a proper subset, *i.e.*, *coreset*, of the training samples. To this end, we run K-means clustering [44] on the normalized features $W_m := z(x_m)/\|z(x_m)\|$ using cosine similarity as a metric. Then, we use the center of each cluster as the coreset. For contrasting shifted instances (4), we choose the coreset for each shifting transformation S . Table 12 shows the results for various coreset sizes, given by a ratio from the full training samples. Keeping only a few (*e.g.*, 1%) samples is sufficient.

Table 12: AUROC (%) values of our method for various coreset sizes (% of training samples), under one-class (OC-) CIFAR-10, CIFAR-100 (super-class), and ImageNet-30. The values are averaged over classes. Keeping only a few (*e.g.*, 1%) samples shows sufficiently good results.

Coreset (%)	OC-CIFAR-10	OC-CIFAR-100	OC-ImageNet-30
1%	94.22	89.27	91.06
10%	94.30	89.46	91.51
100%	94.31	89.55	91.63

F Ablation study on the balancing terms

We study the effects of the balancing terms λ_S^{con} , λ_S^{cls} in Section 2.3. To this end, we compare of our final loss (5), without (w/o) and with (w/) the balancing terms λ_S^{con} and λ_S^{cls} . When not using the balancing terms, we set $\lambda_S^{\text{con}} = \lambda_S^{\text{cls}} = 1$ for all S . We follow the experimental setup of Table 1, *e.g.*, use rotation for the shifting transformation. We run our experiments on CIFAR-10, CIFAR-100 (super-class), and ImageNet-30 datasets. Table 13 shows that the balancing terms gives a consistent improvement. CIFAR-10 do not show much gain since all λ_S^{con} and λ_S^{cls} show similar values; in contrast, CIFAR-100 (super-class) and ImageNet-30 show large gain since they varies much.

Table 13: AUROC (%) values of our method without (w/o) and with (w/) balancing terms, under one-class (OC-) CIFAR-10, CIFAR-100 (super-class), and ImageNet-30. The values are averaged over classes, and bold denotes the best results. Balancing terms give consistent improvements.

	OC-CIFAR-10	OC-CIFAR-100	OC-ImageNet-30
CSI (w/o balancing)	94.28	89.00	91.04
CSI (w/ balancing)	94.31	89.55	91.63

G Combining multiple shifting transformations

We find that combining multiple shifting transformations: given two transformations \mathcal{S}_1 and \mathcal{S}_2 , use $\mathcal{S}_1 \times \mathcal{S}_2$ as the combined shifting transformation, can give further improvements. Table 14 shows that combining “Noise”, “Blur”, and “Perm” to “Rotate” gives additional gain. We remark that one can investigate the better combination; we choose rotation for our experiments due to its simplicity.

Table 14: AUROC (%) values of our method under various shifting transformations. Combining “Noise”, “Blur”, and “Perm” to “Rotate” gives additional gain.

	Base	Noise	Blur	Perm	Rotate	Rotate+Noise	Rotate+Blur	Rotate+Perm
AUROC	87.89	89.29	89.15	90.68	94.31	94.65	94.66	94.60

H Discussion on the features of the contrastive score (6)

We find that the two features: a) the *cosine similarity* to the nearest training sample in $\{x_m\}$, *i.e.*, $\max_m \text{sim}(z(x_m), z(x))$, and (b) the *feature norm* of the representation, *i.e.*, $\|z(x)\|$, are important features for detecting OOD samples under the SimCLR representation.

In this section, we first demonstrate the properties of the two features under vanilla SimCLR. While we use the vanilla SimCLR to validate they are general properties of SimCLR, we remark that our training scheme (see Section 2.2) further improves the discrimination power of the features. Next, we verify that cosine similarity and feature norm are *complementary*, that combining both features (*i.e.*, s_{con} (6)) give additional gain. For the latter one, we use our final training loss to match the reported values in prior experiments, but we note that the trend is consistent among the models.

First, we demonstrate the effect of cosine similarity for OOD detection. To this end, we train vanilla SimCLR using CIFAR-10 and CIFAR-100 and in- and out-of-distribution datasets. Since SimCLR attracts the same image with different augmentations, it learns to cluster similar images; hence, it shows good discrimination performance measured by linear evaluation [5]. Figure 3a presents the t-SNE [43] plot of the normalized features that each color denote different class. Even though SimCLR is trained in an unsupervised manner, the samples of the same classes are gathered.

Figure 3b and Figure 3c presents the histogram of the cosine similarities from the nearest training sample (*i.e.*, $\max_m \text{sim}(z(x_m), z(x))$), for training and test datasets, respectively. For the training set, we choose the second nearest sample since the nearest one is itself. One can see that training samples are concentrated, even though contrastive learning pushes the different samples. It complements the results of Figure 3a. For test sets, the in-distribution samples show a similar trend with the training samples. However, the OOD samples are farther from the training samples, which implies that the cosine similarity is an effective feature to detect OOD samples.

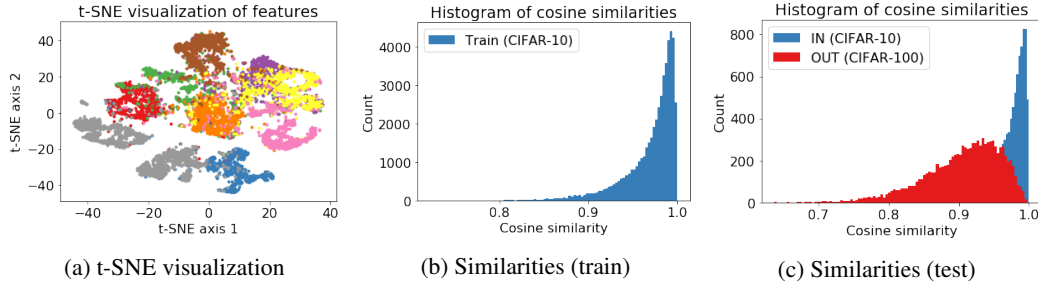


Figure 3: Plots for cosine similarity.

Second, we demonstrate that the feature norm is a discriminative feature for OOD detection. Following the prior setting, we use CIFAR-10 and CIFAR-100 for in- and out-of-distribution datasets, respectively. Figure 4a shows that the discriminative power of feature norm improves as the training epoch increases. We observe that this phenomenon consistently happens over models and settings; the contrastive loss makes the norm of in-distribution samples relatively larger than OOD samples. Figure 4b shows the norm of CIFAR-10 is indeed larger than CIFAR-100, under the final model.

This is somewhat unintuitive since the SimCLR uses the *normalized* features to compute the loss (1). To understand this phenomenon, we visualize the t-SNE [43] plot of the feature space in Figure 4c, randomly choosing 100 images from both datasets. We randomly augment each image for 100 times for better visualization. One can see that in-distribution samples tend to be spread out over the large sphere, while OOD samples are gathered near center.⁴ Also, note that the same image with different augmentations are highly clustered, while in-distribution samples are slightly more assembled.⁵

We suspect that increasing the norm may be an *easier* way to maximize cosine similarity between two vectors: instead of directly reducing the feature distance of two augmented samples, one can also increase the overall norm of the features to reduce the *relative* distance of two samples.

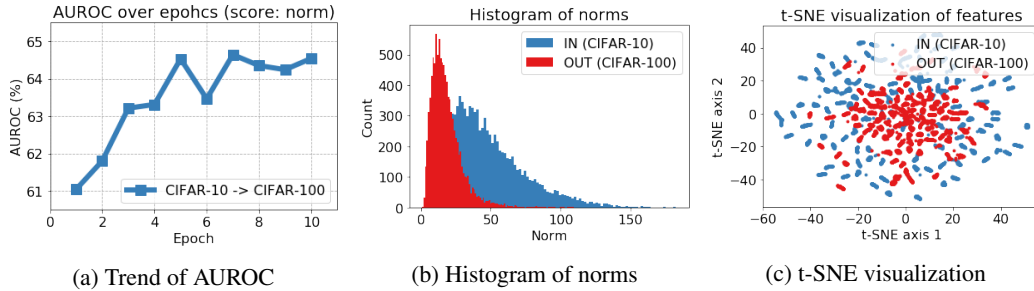


Figure 4: Plots for feature norm.

Finally, we verify that cosine similarity (sim-only) and feature norm (norm-only) are complementary: combining them (sim+norm) gives additional improvements. Here, we use the model trained by our final objective (5), and follow the inference scheme of the main experiments (see Table 7). Table 15 shows AUROC values under sim-only, norm-only, and sim+norm scores. Using only sim or norm already shows good results, but combining them shows the best results.

Table 15: AUROC (%) values for sim-only, norm-only, and sim+norm (*i.e.*, contrastive (6)) scores, under one-class (OC-) CIFAR-10, CIFAR-100 (super-class), and ImageNet-30. The values are averaged over classes. Using both sim and norm features shows the best results.

	OC-CIFAR-10	OC-CIFAR-100	OC-ImageNet-30
Sim-only	90.12	86.57	83.18
Norm-only	92.70	87.71	88.56
Sim+Norm	93.32	88.79	89.32

⁴t-SNE plot *does not* tell the true behavior of the original feature space, but it may give some intuition.

⁵We also try the local variance of the norm as a detection score. It also works well, but the norm is better.

I Rethinking OOD detection benchmarks

We find that resized LSUN and ImageNet [39], one of the most popular benchmark datasets for OOD detection, are visually far from in-distribution datasets (commonly, CIFAR [33]). Figure 5 shows that resized LSUN and ImageNet contain artificial noises, produced by broken image operations.⁶ It is problematic since one can detect such datasets with simple data statistics, without understanding semantics from neural networks. To progress OOD detection research one step further, one needs more *hard* or *semantic* OOD samples that cannot be easily detected by data statistics.

To verify this, we propose a simple detection score that measures the *input smoothness* of an image. Intuitively, noisy images would have a higher variation in input space than natural images. Formally, let $x^{(i,j)}$ be the i -th value of the vectorized image $x \in \mathbb{R}^{HWK}$. Here, we define the *neighborhood* \mathcal{N} as the set of spatially connected pairs of pixel indices. Then, the *total variation* distance is given by

$$\text{TV}(x) = \sum_{i,j \in \mathcal{N}} \|x^{(i)} - x^{(j)}\|_2^2. \quad (14)$$

Then, we define the *smoothness score* as the difference of total variation from the training samples:

$$s_{\text{smooth}}(x) := |\text{TV}(x) - \frac{1}{M} \sum_m \text{TV}(x_m)|. \quad (15)$$

Table 16 shows that this simple score detects current benchmark datasets surprisingly well.

To address this issue, we construct new benchmark datasets, using a fixed resize operation⁷, hence coined LSUN (FIX) and ImageNet (FIX). For LSUN (FIX), we randomly sample 1,000 images from every ten classes of the training set of LSUN. For ImageNet (FIX), we randomly sample 10,000 images from the entire training set of ImageNet-30, excluding “airliner”, “ambulance”, “parking-meter”, and “schooner” classes to avoid overlapping with CIFAR-10.⁸ Figure 6 shows that the new datasets are more visually realistic than the former ones (Figure 5). Also, Table 16 shows that the fixed datasets are not detected by the simple data statistics (15). We believe our newly produced datasets would be a stronger benchmark for hard or semantic OOD detection for future researches.

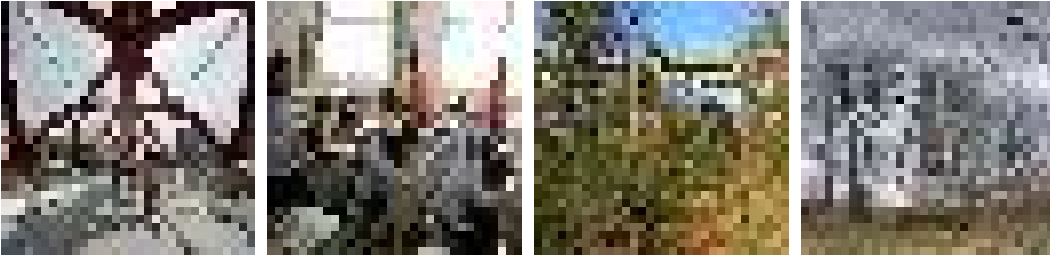


Figure 5: Current benchmark datasets: resized LSUN (left two) and ImageNet (right two).

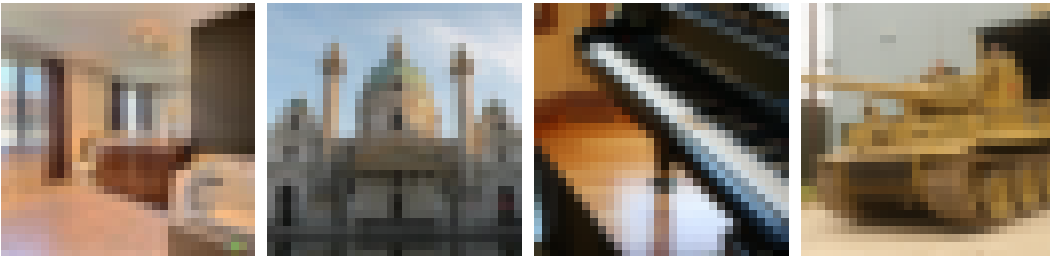


Figure 6: Proposed datasets: LSUN (FIX) (left two) and ImageNet (FIX) (right two).

⁶It is also reported in <https://twitter.com/jaakkolehtinen/status/1258102168176951299>.

⁷We use PyTorch `torchvision.transforms.Resize()` operation.

⁸We provide the datasets and data generation code in <https://github.com/alinelab/CSI>.

Table 16: AUROC (%) values using the smoothness score (15), under unlabeled CIFAR-10. Bold denotes the values over 80%, which implies the dataset is easily detected.

CIFAR10 →						
SVHN	LSUN	ImageNet	LSUN (FIX)	ImageNet (FIX)	CIFAR-100	Interp.
85.88	95.70	90.53	44.13	52.76	52.14	66.17

J Additional examples of rotation-invariant images

We provide additional examples of rotation-invariant images (see Table 6 in Section 3.2). Those image commonly appear in real-world scenarios since many practical applications deal with non-natural images, *e.g.*, manufacturing - steel [62] or textile [60] for instance, or aerial [70] images. Figure 7 and Figure 8 visualizes the samples of manufacturing and aerial images, respectively.

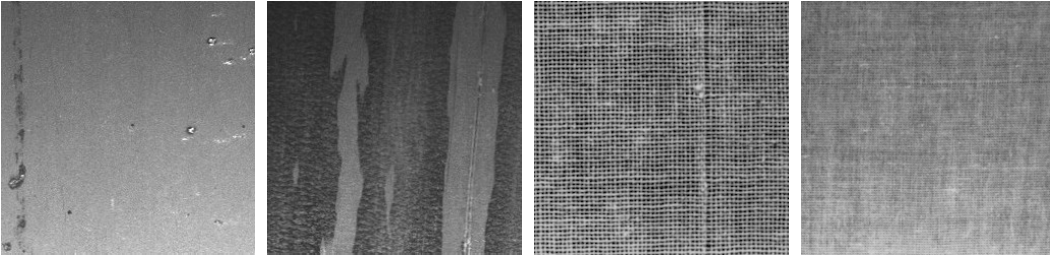


Figure 7: Examples of steel (left two) and textile (right two) images.



Figure 8: Examples of aerial images.