Three approaches to facilitate DNN generalization to objects in out-of-distribution orientations and illuminations: late-stopping, tuning batch normalization and invariance loss

Akira Sakai^{a,*}, Taro Sunagawa^a, Spandan Madan^{b,d}, Kanata Suzuki^a, Takashi Katoh^a, Hiromichi Kobashi^a, Hanspeter Pfister^b, Pawan Sinha^c, Xavier Boix^{c,d,1,*}, Tomotake Sasaki^{a,d,1,*}

^aFujitsu Limited, Kawasaki, Japan. ^bHarvard University, Cambridge, USA. ^cMassachusetts Institute of Technology, Cambridge, USA. ^dCenter for Brains, Minds and Machines, Cambridge, USA.

Abstract

The training data distribution is often biased towards objects in certain orientations and illumination conditions. While humans have a remarkable capability of recognizing objects in out-of-distribution (OoD) orientations and illuminations, Deep Neural Networks (DNNs) severely suffer in this case, even when large amounts of training examples are available. In this paper, we investigate three different approaches to improve DNNs in recognizing objects in OoD orientations and illuminations. Namely, these are (i) training much longer after convergence of the in-distribution (InD) validation accuracy, *i.e.*, late-stopping, (ii) tuning the momentum parameter of the batch normalization layers, and (iii) enforcing invariance of the neural activity in an intermediate layer to orientation and illumination conditions. Each of these approaches substantially improves the DNN's OoD accuracy (more than 20% in some cases). We report results in four datasets: two datasets are modified from the MNIST and iLab datasets, and the other two are novel (one of 3D rendered cars and another of objects taken from various controlled orientations and illumination conditions). These datasets allow to study the effects of different amounts of bias and are challenging as DNNs perform poorly in OoD conditions. Finally, we demonstrate that even though the three approaches focus on different aspects of DNNs, they all tend to lead to the same underlying neural mechanism to enable OoD accuracy gains—individual neurons in the intermediate layers become more selective to a category and also invariant to OoD orientations and illuminations.

Keywords: Out-of-distribution Generalization; Object Recognition in Novel Illuminations and Orientations; Neural Selectivity and Invariance

1. Introduction

The object recognition performance of Deep Neural Networks (DNNs) dramatically degrades when the train and test distributions are not identical due to dataset bias [1], *i.e.*, when tested in out-of-distribution (OoD) conditions. There is a big gap between DNNs and humans when evaluated in OoD conditions. This issue has been getting much interest in recent years [2, 3, 4, 5, 6], as it severely compromises the safety and fairness of AI applications.

One of the most prominent factors of dataset bias is that objects may appear in a constrained range of orientation and illumination conditions [7, 8]. While generalization to OoD orientations and illumination conditions has been long studied in both biological and artificial neural networks, *e.g.*, [9, 10, 11], the computational mechanisms that facilitate such generalization remain as a key outstanding question. Recently, [12, 13] have shown that DNNs are capable to overcome bias by transferring the generalization ability obtained from objects seen in a richer set of conditions to the objects seen in biased conditions. Also, the emergence of representations at the individual neuron level in the intermediate layers of the DNN that are selective to categories and invariant to the OoD conditions has been identified as

^{*}Corresponding author

Email addresses: akira.sakai@fujitsu.com(Akira Sakai), xboix@mit.edu(Xavier Boix),

tomotake.sasaki@fujitsu.com(Tomotake Sasaki) ¹Equal contribution

a mechanism that may facilitate such OoD generalization. Invariant neural representations have been studied during decades, *e.g.*, [11], and here they appear as the mechanism that allows OoD generalization. This begs the question whether we can further encourage the emergence of invariant neural representations in DNNs in order to further improve OoD generalization.

In this paper, we investigate factors that can substantially boost the DNN ability to recognize objects in OoD orientations and illuminations. In particular, we discover that the following factors, summarized in Fig. 1, have a remarkable impact:

- Late-stopping: DNNs are usually trained until the validation recognition accuracy (which is indistribution) converges. We found that in many cases the OoD recognition accuracy improves slowly, yet consistently, after the validation (indistribution) accuracy has converged. This finding is surprising as classic machine learning theory suggests early-stopping as a regularization mechanism [14], and we found that the opposite is beneficial to improve OoD generalization in DNNs. We call this approach "late-stopping".
- 2. *Tuning the batch normalization parameter:* Batch normalization (BN) is known to have an impact in OoD recognition accuracy [15]. We found that tuning the only hyperparameter of BN, *i.e.*, the momentum, yields substantial gains of OoD recognition accuracy. This approach is denoted as "tuned BN".
- 3. Neural activity invariance loss: Motivated by the aforementioned finding in previous works that invariant neural representations leads to improvements of the OoD recognition accuracy, we include an additional term in the loss function to encourage this phenomenon. This loss term takes the Euclidean distance between neural activity corresponding to pairs of images from the same category on an intermediate layer. By minimizing this loss term, the neural activity tends to be invariant for objects of the same category even in different viewing conditions. We do not consider that pairs of images from different categories should have distinguishable neural activity, since the classification loss term already encourages this. We call this approach "invariance loss" in short.

Our results demonstrate that each of these three approaches alone lead to substantial improvements of object recognition in OoD orientations and illumination conditions. Results also corroborate that when any of the three approaches leads to an increase of selectivity and invariance at the individual neuron level, OoD recognition accuracy improves in the majority of trials. Experiments are performed in four challenging benchmarks, namely modifications of the MNIST dataset [16] and iLab dataset [17] and two novel datasets we introduce, which are the CarsCG and the MiscGoods datasets. CarsCG contains 3D rendered cars from different orientations, and the MiscGoods dataset consists of images of objects taken with a robotic arm from different viewpoints and controlled illumination conditions. These datasets allow to evaluate the DNN generalization ability to recognize objects in OoD orientations and illumination conditions. Also, they allow to analyze the effects of different amounts of bias and are challenging as DNNs perform poorly in OoD conditions.

2. Previous works

Our results add to the growing body of literature to improve the generalization ability of DNNs to OoD orientations and illumination conditions. Prior efforts leverage synthesized sources of training data [18, 19, 20, 21], 3D models of objects [22], specific characteristics of the target domain [23, 24, 25], or sensing approaches such as omnidirectional imaging [26]. These approaches add preconceived components to the DNN that need to be adjusted at hand for new objects and conditions. Here, we focus on pure learning-based strategies as these are not constrained to specific objects and conditions and can be automatically adjusted to new datasets.

Other strands of research that live in neighbouring areas investigate generalization to new domains and also, overcoming spurious correlations between image features and categories. While domain generalization does not tackle dataset bias and overcoming spurious correlations does not address recognition of objects in OoD orientations and illuminations, these two research areas use related techniques and concepts to our work. In the following we review both of them.

Domain generalization. There is a plethora of works that consists on learning representations in several domains that can be easily transferred to new domains, *e.g.*, [27, 28, 29, 30, 31, 32, 33, 34]. The problem of domain generalization is similar to the problem overcoming dataset bias in our study in the sense that representations that facilitate generalization to novel conditions should be learned. However, in domain generalization the learner has access to multiple domains during training that can be leveraged for generalization,



Figure 1: *Three approaches to facilitate generalization to objects in out-of-distribution (OoD) orientations and illuminations.* (a) Learning curves of in-distribution (InD) test accuracy and OoD accuracy for late-stopping applied to the MiscGoods-illuminations dataset (medium InD data diversity). OoD accuracy converges much later than InD accuracy. (b) Learning curves of the OoD accuracy with and without tuning batch normalization momentum (tuned BN) in the CarsCG-Orientations, dataset (medium InD data diversity). It can be seen that tuning the momentum reduces the oscillation of the OoD accuracy and improves the performance. (c) Left: Conceptual diagram of the invariance loss. Pairs of images that belong to the same category are fed into the DNN. The invariance loss is based on the Euclidean distance between the pairs of the last ReLU activity. The classification loss is calculated with the network output as usual. The total loss is the weighted sum of the invariance and classification losses. Right: Learning curve of OoD accuracy in MiscGoods-illuminations dataset (medium InD data diversity) when the invariance loss is applied. The OoD accuracy increases by about 20% compared to the baseline. The solid lines in the plots are the mean value. The lighter semitransparent colors surrounding the the solid lines indicate 95% confidence interval.



Figure 2: *Sample images of four datasets*. (a) MNIST-Positions, (b) iLab-Orientations, (c) CarsCG-Orientations, and (d) MiscGoods-Illuminations are shown in each subfigure. Samples from each dataset are arranged in a grid pattern. Each row indicates categories and each column indicates either an orientation or an illumination condition.

while in the problem of overcoming dataset bias only one training set is available. Recently, several works in domain generalization [35, 36] highlighted the need of invariant representations to obtain further improvements in generalization, which further motivates investigating the invariance loss in our study.

Overcoming spurious correlations between image features and categories. Many datasets are biased in a way that a specific image feature consistently appears in images of the same category. DNNs tend to learn that those features are informative of the category [37]. This form of dataset bias is different from the bias in the object orientation and illumination conditions, which do not necessarily lead to spurious correlations. Recently, there have been a several works that address spurious correlations. These are based on automatically detecting the features that spuriously correlate with the category, and encourage the DNN not to rely on those features [38, 39]. Ahmed et al. [40] introduced a method that effectively alleviates the effect of spurious correlation caused by biased object background. This work exploits the assumption that the training distribution also contains examples without spurious correlations. It employs EIIL [41] to classify the images of an category with the features that spuriously correlate with the category and without them. Then, invariance is encouraged across these two groups of images. Thus, invariance appears once more as a facilitator of generalization.

3. Performance degradation on OoD Conditions

In this section, we introduce the methodology to evaluate the accuracy of the DNN in OoD conditions. First, we describe the procedure of the bias-controlled experiment. Next, we introduce the four datasets used in this study and finally, we evaluate the performance degradation that occurs in OoD conditions in these four datasets.

3.1. Bias-controlled experiments

In a dataset there could be multiple biasing factors at the same time that can cause performance degradation. In the datasets in this study, we analyze either the orientation or illumination condition, as it allows to more clearly understand the effect of each individual factor. Thus, the datasets we use contains several



Figure 3: InD and OoD combinations for biascontrolled experiments. Each sample is a combination of a category and an orientation or illumination condition. We create a set of combinations called "InD combinations" and a set of combinations called "OoD combinations". The ratio of InD combinations to all combinations is called InD data diversity. In addition, we create a train dataset $(\mathcal{D}_{train}^{(InD)})$ and an InD validation dataset $(\mathcal{D}_{val}^{(InD)})$ from samples included in the InD combinations, and an OoD test dataset $(\mathcal{D}^{(OoD)})$ from the samples included in the OoD combinations.

combinations of categories and orientation or illuminations conditions. We use C to denote the set of all categories and N the set of all orientation or illuminations conditions. Let $\mathbf{x}^{(k)}$ be an image of the dataset and let $\mathbf{y}^{(k)} := (c^{(k)}, n^{(k)})$ be a tuple representing the groundtruth category (*i.e.*, $c^{(k)} \in C$), and the orientation or illuminations condition (*i.e.*, $n^{(k)} \in N$).

In order to evaluate the DNN's OoD generalization capabilities, we train them in a dataset that follows a distribution that only contains a subset of all possible combinations, *i.e.*, a subset of $C \times N$. Then, the DNN is evaluated with images from combinations that were not included in the training distribution. Let $\mathcal{I} \subset C \times N$ be the set of combinations used to generate the InD combinations. We ensure that \mathcal{I} contains all categories and all conditions at least once (but not all combinations), such that we have images from all image categories and conditions in a balanced manner.

We use $\mathcal{D}^{(\text{InD})}$ to denote the set of images that are InD, *i.e.*, images whose label is in \mathcal{I} , $\mathbf{y}^{(k)} \in \mathcal{I}$. Namely, the InD images dataset, $\mathcal{D}^{(\text{InD})}$, is defined as in the following:

$$\mathcal{D}^{(\mathrm{InD})} := \{ (\mathbf{x}, \mathbf{y}) | \mathbf{y} \in \mathcal{I} \}.$$
(1)

 $\mathcal{D}^{(\mathrm{InD})}$ is further divided into train dataset and valida-

tion dataset, which we denote as $\mathcal{D}_{train}^{(InD)}$ and $\mathcal{D}_{val}^{(InD)}$, respectively. The term InD accuracy refers to the DNN's accuracy on $\mathcal{D}_{val}^{(InD)}$. The OoD dataset $\mathcal{D}^{(OoD)}$ is defined as

$$\mathcal{D}^{(\text{OoD})} := \{ (\mathbf{x}, \mathbf{y}) | \mathbf{y} \in (\mathcal{C} \times \mathcal{N}) \setminus \mathcal{I} \}.$$
(2)

The term OoD accuracy refers to the accuracy on the OoD dataset $\mathcal{D}^{\rm (OoD)}.$

We also define the InD data diversity of a dataset as $\#(\mathcal{I})/\#(\mathcal{C} \times \mathcal{N})$, where $\#(\cdot)$ denotes a number of elements. Thus, the data diversity measures the portion of combinations included in the training distribution. To directly compare the effect of the InD data diversity on the OoD accuracy, we vary the InD data diversity such that the combinations in the distributions of lower InD data diversity are included in the combinations of higher InD data diversity, while keeping the training set size constant, *i.e.*, $\#(\mathcal{D}_{\text{train}}^{(\text{InD})}(\mathcal{I}))$ is constant for all InD data diversity. These restrictions allows us to evaluate the performance of the DNN only by the difference in InD data diversity, not by the difference in the amount of combinations or training examples.

3.2. Datasets

We use the following four datasets. These datasets have labels for both category and either orientation or illumination condition, in order to evaluate OoD generalization. See Appendix A for further details than the ones provided in the following.

MNIST-Positions. It is based on the MNIST dataset [16]. We created a dataset of 42×42 pixels with nine numbers by resizing images to 14×14 and placing them in one of nine possible positions in a 3×3 empty grid. We call this dataset the *MNIST-Positions* dataset. In our experiments, the digits are considered to be the category set, and the positions where the digits are placed is considered as the orientation. We use nine digits and nine positions. Samples are shown in Fig. 2(a). We used 54K images for $\mathcal{D}_{train}^{(InD)}$, 8K images for $\mathcal{D}_{val}^{(InD)}$ and 8K images for $\mathcal{D}_{val}^{(InD)}$. Low, medium, and high InD data diversity are set to be 2/9, 4/9, and 8/9, respectively.

iLab-Orientations. iLab-2M is a dataset created from iLab-20M dataset [17]. The dataset consists of images of 15 categories of physical toy vehicles photographed in various orientations, elevations, lighting conditions, camera focus settings and backgrounds. The image size is 256×256 pixels. From the original iLab-2M dataset, we chose six categories (bus, car, helicopter, monster truck, plane, and tank) and six orientations. We call it



Figure 4: *Performance degradation in the OoD conditions*. Upper figures show the examples of InD combination in MiscGoods-illuminations dataset to depict the InD data diversity of each experiment. Lower figures show InD or OoD accuracy of ResNet-18 in (a) low InD data diversity, (b) medium InD data diversity and (c) high InD data diversity performed on four datasets. Each experiment is conducted five times, and the mean and 95% confidence interval are reported. Sharp performance degradation in OoD accuracy is observed (*e.g.*, between 40% to almost 80% is observed when the InD data diversity is low). These result shows the impact of a distribution shift from InD to OoD to the performance of a DNN.

iLab-Orientations. Samples are shown in Fig. 2(b). We resized each image to 64×64 pixels. We used 18K images for $\mathcal{D}_{train}^{(InD)}$, 8K images for $\mathcal{D}_{val}^{(InD)}$ and 8K images for $\mathcal{D}^{(OoD)}$. Low, medium, and high InD data diversity are set to be 2/6, 3/6, and 5/6, respectively.

CarsCG-Orientations. CarsCG-Orientations is a new dataset that consists of images of ten types of cars in various conditions rendered by Unreal Engine. It includes ten orientations, three elevations, ten body colors, five locations and three time frames (morning, evening, night). We synthesize images with 1920×1080 pixels and resize them as 224×224 pixels for our experiment. We chose ten types of cars as categories and ten orientations for each of them. Samples are shown in Fig. 2(c). More samples are provided in Appendix A. In the experiment, we used 3400 images for $\mathcal{D}_{train}^{(InD)}$, 450 images for $\mathcal{D}_{val}^{(InD)}$ and 800 images for $\mathcal{D}_{cod}^{(OoD)}$. Low, medium, and high InD data diversity are set to be 2/10, 5/10, and 9/10, respectively.

MiscGoods-Illuminations. MiscGoods-Illuminations is a subset of DAISO-10, a novel dataset collected

for this study. The dataset consists of ten physical miscellaneous goods photographed using a robotic arm with five controlled illumination conditions, two object placements, twenty object orientations, and five camera angles. Each image is 640×480 pixels in size. We chose five categories (stuffed dolphin, stuffed whale, metal basket, imitation plant and cup) and five illumination conditions as shown in Fig. 2(d). More samples are displayed in Appendix A. We resize the images to 224×224 pixels for our experiments. We used 800 images for train $\mathcal{D}_{train}^{(InD)}$, 200 images for $\mathcal{D}_{val}^{(InD)}$ and 400 images for $\mathcal{D}_{val}^{(OoD)}$. Low, medium, and high InD data diversity are set to be 2/5, 3/5 and 4/5, respectively.

3.3. OoD accuracy results

We now demonstrate that these four datasets are extremely challenging for DNNs as these achieve low accuracy in OoD conditions. We examine the performance degradation in three InD data diversity: low, medium, and high. Recall that we evaluate InD accuracy in $\mathcal{D}_{val}^{(InD)}$ and the OoD accuracy in $\mathcal{D}^{(OoD)}$. We

use ResNet-18 [42] trained with $\mathcal{D}_{train}^{(InD)}$. The experimental setup is introduced in Sec. 6.1.

Figure 4 shows the OoD accuracy degradation regarding the four datasets ranging low to high InD data diversity. While the InD accuracy is more than 80%for all four datasets at almost all data diversities (except for MNIST-positions), the OoD accuracy showed a substantial degradation when the DNN was trained with low and medium InD data diversities. Between 20% to 70% performance degradation is observed in low InD data diversity in all four datasets. In medium InD data diversity, large performance degradation ranging from 10% to 50% is observed, and for high InD data diversity, there is more than 10% performance degradation in CarsCG-Orientations and MiscGoods-Illuminations datasets. Thus, dramatic drops of accuracy are observed in OoD conditions, which confirms that these benchmarks are very challenging for DNNs.

OoD accuracy is often overlooked in standard computer vision benchmarks and only InD is usually reported. This is usually due to the difficulty of measuring OoD accuracy. Our datasets enable evaluating OoD accuracy in a controlled way that facilitates understanding the different factors that may affect the OoD accuracy. The performance degradation in OoD conditions is expected when deploying application of deep learning. Recently, it has been reported that even a small amount of data bias can cause major performance degradation [5], and this is reconfirmed for our four datasets. Also, the drop of accuracy in our datasets is dramatic, specially for low InD data diversity. Our datasets allow to gain an understanding of the specific biasing factors in the dataset, i.e., orientation and illumination conditions, and analyze aspects such as the InD data diversity.

4. Three approaches to improve OoD Accuracy

We now introduce the three approaches to address the performance drop of accuracy in OoD conditions, which are "late-stopping", "tuning the batch normalization momentum" and "invariance los". These three approaches are independent on each other and tackle different aspects of the DNN training.

4.1. Late-stopping

The stopping criteria for training is known to have an impact on the DNNs performance [43, 44, 45]. In particular, stopping the training before convergence of the training accuracy, *i.e.*, early stopping, is known to prevent overfitting in shallow classifiers [46]. However, these results are with respect to InD accuracy, and little is known regarding the relation between the stopping critaria and OoD accuracy. We therefore run experiments with a large number of training epochs (up to 1000 epochs) in order to investigate any patterns. Figure 1(a) shows the change of InD and OoD accuracy when ResNet-18 is trained with the medium InD data diversity. Surprisingly, the OoD accuracy, unlike the InD accuracy, continued to increase in performance after training during a large number of epochs. This phenomenon was not known because usually only the InD accuracy is analyzed. We denote the approach of continuing the training of a DNN after the convergence of InD validation accuracy as "late stopping".

4.2. Tuning batch normalization

Batch normalization [47] is a method used to speedup and stabilize the training of DNN networks through normalization of the layers' inputs by re-centering and re-scaling them. Batch normalization has also been reported to act as a regularizer and improve generalization [48]. Thus, it is reasonable that batch normalization could help improving OoD generalization but this has not been studied so far.

Batch normalization uses the so called moving average to recenter the layer's input. Let $v_{ma}(t)$ be the moving average at training step t. The moving average is updated at each training step in the following way:

$$\boldsymbol{v}_{\mathrm{ma}}(t) = (\beta - 1)\boldsymbol{v}_{\mathrm{mean}}(t) + \beta \boldsymbol{v}_{\mathrm{ma}}(t-1), \quad (3)$$

where $v_{\text{mean}}(t)$ is the mean activity over the batch of the *t*-th training step, and $\beta \in [0, 1]$ is called momentum and balances the update of the moving average between $v_{\text{mean}}(t)$ and itself. Note that the only hyperparameter available for batch normalization is β , and we use this to adjust it. Usually, β is set to 0.9 or 0.99, which is the default value. We use the default value 0.99 that is employed by the TensorFlow library [49].

We investigated how the OoD generalization performance behaves depending on the value of the batch normalization momentum, β . Figure 1(b) shows the learning curves of ResNet-18 trained on MiscGoodsilluminations with the medium InD data diversity. Experimentally, we found that the tuning momentum parameter, β , can have a significant positive impact on the OoD generalization performance. Generally, the default value of $\beta = 0.99$ was too large for almost all cases in our experiments. We call this approach as tuning batch normalization or "tuning BN".

4.3. Invariance loss

The "invariance loss" approach is intended to increase the invariance score that is introduced in Madan *et al.* [13], which we explain in Sec. 5. This invariance score measures the degree of invariance in the neural activity of intermediate layers, and previous works have shown that DNNs that generalize better to OoD conditions have developed larger degrees of invariance in the intermediate layers.

Concretely, we encourage the emergence of invariant representations by taking pairs of images that belong to the same category and enforce that the neural activity is as similar as possible. To do so, we use the Euclidean distance between the activities of neurons in an intermediate layer caused by the pairs of images, and add this as an additional loss term to the classification loss. Figure 1(c) shows the scheme of this approach. Let $g(\cdot; \theta_a)$ be the neural activity of a DNN intermediate layer, where θ_a are the parameters of the DNN before the intermediate layer. Let $f(\cdot; \theta_f)$ be the output of the DNN given as input the intermediate layer, $g(\cdot; \theta_q)$, where θ_f are the DNN parameters from the intermediate layer to the output of the network. Thus, the neural activity of the intermediate layer for an image x is $g(x; \theta_q)$ and the output of the whole network is $f(g(\mathbf{x}; \boldsymbol{\theta}_q); \boldsymbol{\theta}_f)$. Let \mathbf{x} be a training image, and let \mathbf{x}' be another image that belongs to the same category as x, and is sampled from the training data $D_{\rm train}^{({\rm InD})}$ according to some sampling strategy (in our experiments is random with uniform distribution across the training images of the same category). Thus, the invariance loss is expressed as

$$\|\boldsymbol{g}(\mathbf{x};\boldsymbol{\theta}_g) - \boldsymbol{g}(\mathbf{x}';\boldsymbol{\theta}_g)\|_2.$$
(4)

This term is added to the categorical cross entropy loss weighted with a hyperparameter that we call λ , such that the invariance loss term acts as a regularization term. Note that the invariance loss is equivalent to the contrastive loss [50] for positive examples in the context of metric learning, but it has not been used so far to improve generalization to OoD orientations and illumination conditions.

5. Selectivity and invarinace for OoD generalization

We now revisit the mechanism at the individual neuron level of intermediate layers that previous works have suggested that facilitates OoD generalization, *i.e.*, individual neurons being selective to a category and invariant OoD conditions. This mechanism has been shown to explain the improvement in OoD accuracy with increased InD data diversity [12, 13].

For a given intermediate layer of the DNN, let α_{cn}^j be the average activity for the *j*-th neuron over all images with the *c*-th category and the *n*-th orientation or illumination condition. For neuron *j*, the activity is 0-1 normalized. Let c^{*j} be the category that a neuron *j* is most active on average, *i.e.*, $c^{*j} := \operatorname{argmax}_c \sum_n \alpha_{cn}^j$. This is called preferred category. The selectivity score S^j is defined as

$$S^{j} := \frac{\hat{\alpha}^{j} - \bar{\alpha}^{j}}{\hat{\alpha}^{j} + \bar{\alpha}^{j}},\tag{5}$$

where, $\hat{\alpha}^{j} := \frac{1}{\#(\mathcal{N})} \sum_{n} \alpha_{c^{*j}n}^{j}$ and $\bar{\alpha}^{j} := \sum_{n \in \mathcal{N}} \alpha_{c^{*j}n}^{j}$

 $\frac{\sum_{c \neq c^* d} \sum_N \alpha_{cn}^j}{\#(C)(\#(C)-1)}$ denote the average activity for the preferred category and for the remaining categories, respectively. This selectivity score ranges from zero to one and takes its maximum value in the case that the neuron average activity, α_{cn}^j , is 0 for all categories except for the preferred category, *i.e.*, the neuron is only active for the preferred category. The invariance score I^j is defined as

$$I^{j} := 1 - (\max_{n} \alpha_{c^{*j}n}^{j} - \min_{n} \alpha_{c^{*j}n}^{j}), \qquad (6)$$

and it also ranges from zero to one and takes the maximum in the case that the average activity, α_{cn}^{j} , takes the same value for the preferred category regardless of the orientation and illumination conditions.

Finally, we define the SI score of a neuron as the geometric mean of the selectivity and invariance scores, *i.e.*, $\sqrt{S^{j}I^{j}}$. Neurons that have a larger SI score are active for specific categories independently on the orientation and illumination conditions. Networks with neurons that have larger SI scores have been observed to generalize better in OoDconditions. In order to provide a score that summarizes the SI score across all neurons in the layer, we use the upper 20 percentile of the scores among all neurons. This is because not all neurons are required to have larger SI to improve OoD generalization, and we just take into account a portion of neurons with the highest SI score. In the experiment section, we use this summary of the SI score across neurons to assess whether the three approaches we introduce yield improved OoD accuracy through improving selectivity and invariance.

6. Experiments

We first introduce the experimental setup, and then report the OoD accuracy facilitated by the three approaches explained in Sec. 4. Finally, we analyze whether this boost of OoD accuracy is driven by selective and invariance mechanism revisited in Sec. 5.

6.1. Experimental setting

We apply the three approaches to improve OoD accuracy to ResNet-18 [42] and evaluate its effectiveness in the aforementioned datasets (MNIST-Positions, iLab-Orientations, CarsCG-Orientations, and MiscGoods-Illuminations). Standard ResNet-18 is adopted as the network for all experiments and we trained it in the standard manner. Namely, all neurons employ the ReLU activation function $g(z) = \max\{0, z\}$ [51] and Glorot uniform initializer [52] is adopted for the network weights initialization for all experiments. Adam [53] is employed as the optimization algorithm. The pixels of images are normalized within 0 to 1 as a preprocessing for all datasets.

We run five trials in all cases and report mean accuracy and its 95% confidence interval. In each trial, the InD combinations are chosen randomly as long as they satisfy the conditions explained in Sec. 3.1, and the OoD combinations are created accordingly. Each of the four approaches, including baseline, is subjected to a hyper-parameter search before performing the five We select the hyper-parameters in a differtrials. ent trial from the ones used to report OoD accuracy. In this reserved trial, we select the hyper-parameters with the highest OoD accuracy by grid search. For all tested approaches, we selected a learning rate in $\{0.1, 0.01, 0.001, 0.0001, 0.00001\}$, and other hyperparameters depending in the approach. In the following we detail the experimental setting of the different approaches.

Late-stopping. The epoch size is set to 1,000 epochs for late stopping, and 100 epochs for the other approaches, including baseline. We confirmed that 100 epochs are sufficient for convergence in InD accuracy by the preliminary experiments. For late stopping, we run as many epochs as computing resources allow (about a week of training).

Tuning batch normalization. For tuning batch normalization, we perform a grid-search for $\beta = \{0.01, 0.1, 0.5, 0.9, 0.99\}$ in addition to the learning rate. For the other approaches, we use 0.99 as a momentum parameter β for batch normalization layer, which is the default value.

Invariance loss. Invariance loss is applied to the last ReLU activation layer "activation_17" which has 512 neurons. We keep fixed the pairs of images in which

invariance is enforced, and we randomize the pairs from time to time. We perform a grid search to determine how frequently we randomize the pairs of images (the choices are randomizing every $\{10, 20, 50, 100\}$ epochs). The weight of the invariance loss term, λ , is also selected via a grid search among the following values: $\lambda = \{1.0, 0.1, 0.01, 0.001, 0.0001\}$.

For more details we refer the reader to Appendix Appendix B.

6.2. Improvement of OoD accuracy

Figure 5 compares of mean OoD accuracy between the baseline and the three approaches for all tested data sets and all tested InD data diversities. Looking at the case of the CarsCG-Orientations and MiscGoodsilluminations datasets, we can see that the three approaches increase the mean OoD accuracy at almost all the data diversities. Comparing the three approaches, late stopping and invariance loss both achieves the best improvement rate in some combinations, and batch norm momentum does not achieve the best improvement in any combination. The highest improvement of 22.2% is achieved by late stopping with a high InD data diversity. The performance improvement across datasets and data diversities is remarkable. Only in iLab-Orientations dataset is relatively small, but for high InD data diversity in this dataset, all three approaches achieve better OoD accuracy than the baseline approach. For MNIST-Positions, all three approaches showed an improvement in performance with medium InD diversity. In Appendix C we report the learning curves and the InD accuracy for a more detailed depiction of the effects of the three approaches during training.

We also investigated whether the three approaches combined together are more effective than the best of three approaches applied alone. Thus, we trained networks using late-stopping, tuned BN and invariance loss together. We call this approach "three approaches together". Another way of combining the three approaches is training networks with each approach alone and then selecting the best of the approaches using a validation set. We call this approach "best of three approaches alone". The hyper-parameter tuning method of these combined approaches is detailed in Appendix D. Table 1 shows the comparison between these two combination approaches and also the baseline, *i.e.*, the network trained without any approach to improve the OoD accuracy. The table reports the number times a method is better than another method across all datasets and InD data diversity. The results show that using the



Figure 5: *Performance improvement of the mean OoD accuracy.* Top figures show the examples of InD combinations, *i.e.*, InD data diversity, in MiscGoods-illuminations dataset. (a), (b), and (c) mean OoD accuracy of the three approaches and the baseline for different InD data diversities. Error bars show 95% confidence interval. (d), (e), and (f) increase of the mean OoD accuracy by the three approaches from the baseline. The unit "pp" in figures denotes percentage points, *i.e.*, the unit for the arithmetic difference of two percentages.

best of the three approaches alone obtains the best results in the vast majority of experiments. Interestingly, the three approaches together performs worse than the baseline for more than half of the experiments. This indicates that the three approaches together interfere with each other and should not be used.

6.3. Analysis for selectivity and invariance mechanism

Figure 6 shows the relationship between the SI score of the last ReLU layer and the OoD accuracy for all combinations of dataset, InD data diversity, and approach (details are provided in Fig. C.11). We can see that there is a large correlation between SI score and OoD accuracy (Pearson's correlation coefficient is 0.891). While it has already been shown in Madan et al. [13] that increasing the InD data diversity improves the OoD accuracy and the SI score, here we show for the first time that approaches that targets improving the OoD accuracy also yield increases of the SI score.

Next, we analyze the relationship between improvements of OoD accuracy and increases of the SI score. We investigate whether increases of the SI score always precede improvement of OoD accuracy, which serves to assess whether invariant representations drive OoD generalization in a more stringent way than the correlational analysis presented before. Let $P(\Delta_{\rm acc}^+)$ be the probability that the OoD accuracy increases when using one of the three approaches to train the network, compared to not using it. Also, let $P(\Delta_{\rm SI}^+)$ be the probability that the SI scores increases when using one of the three approaches, compared to not using it. The conditional probabilities between these two events provides insights regarding whether increases of the the SI score precedes the improvements of the OoD accuracy. We

Table 1: *Comparison of ways to combine the three approaches.* We compare the the best of the three approaches alone (*i.e.*, training a network different times each with one of the three approaches alone, and then selecting the best of the three in a validation set), training with the three approaches together (*i.e.*, training a network using the three approaches together), and the baseline (*i.e.*, training the network without using any approach). Results compare how many times each of these strategies is better than another strategy, across InD data diversities in each of the four datasets.

Comparison of OoD accuracy	MNIST	iLab	CarsCG	MiscGoods	Total
Best of 3 approaches alone vs Baseline	3 vs 0	2 vs 1	3 vs 0	3 vs 0	11 vs 1
3 approaches together vs Baseline	1 vs 2	1 vs 2	1 vs 2	2 vs 1	5 vs 7
Best of 3 approaches alone vs 3 approaches together	3 vs 0	2 vs 1	3 vs 0	1 vs 2	9 vs 3

Table 2: Analysis of the dependency between improvements of OoD accuracy and SI score. This table shows the relative frequency of improvement (+) or degradation (-) of the mean OoD accuracy Δ_{acc} or mean SI score Δ_{SI} . Relative frequency P(x) is calculated by counting the number of cases that satisfy the condition $x \in {\{\Delta_{acc}^+, \Delta_{SI}^+\}}$, and normalize it by total number of cases (*i.e.*, 12, 3 possible InD data diversity × 4 datasets). Conditional relative frequency P(y|x) is also calculated by counting the number combinations satisfying $y \in {\{\Delta_{acc}^+, \Delta_{acc}^-\}}$ in the condition of $x \in {\{\Delta_{SI}^+, \Delta_{SI}^-\}}$, and divide it by the number of combinations satisfying x. The first and second columns show the portion of case the mean OoD accuracy Δ_{acc}^+ and the mean SI score Δ_{SI}^+ increased, respectively. The third column shows the portion of cases the mean OoD accuracy increased Δ_{acc}^+ when the mean SI score increased Δ_{SI}^+ . The fourth column shows the portion cases the mean OoD accuracy increased Δ_{acc}^- when the mean SI score increased Δ_{SI}^- .

Approach	$P(\Delta_{\rm acc}^+)$	$P(\Delta_{\rm SI}^+)$	$P(\Delta_{\rm acc}^+ \Delta_{\rm SI}^+)$	$P(\Delta_{\rm acc}^+ \Delta_{\rm SI}^-)$
Late-Stopping (%)	75.0(9/12)	50.0(6/12)	83.3(5/6)	66.6(4/6)
Tuned BN (%)	75.0(9/12)	83.3 (10/12)	80.0 (8/10)	50.0(1/2)
Invariance Loss (%)	91.7(11/12)	83.3(10/12)	$100.0\ (10/10)$	50.0(1/2)
Total (%)	80.6(29/36)	72.2(26/36)	88.4(23/26)	$60.0 \ (6/10)$

calculate the probabilities by evaluating the frequency that the events happen across datasets and InD data diversities. We report them in Table 2.

We observe by analyzing $P(\Delta_{acc}^+)$ that the OoD accuracy increases very often with the three approaches, at least 75% of the cases. In particular, the OoD accuracy increased 91.7% of the cases for the invariance loss. The analysis of $P(\Delta_{SI}^+)$ shows that tuned BN and invariance loss increase the SI score 83.3% of the cases. This suggests that these two approaches tend to improve the SI score. For late-stopping this trend is not as strong. Yet, when analyzing $P(\Delta_{\rm acc}^+|\Delta_{\rm SI}^+)$, we observe that for the three approaches, increases of the SI score precede the improvements of OoD accuracy (this is in 83.3% (5/6), 80.0% (8/10) and 100% (10/10) of the cases for late-stopping, tuned BN and invariance loss, respectively). Note that the invariance loss directly encourages to increase the SI score, and when the SI score in fact increases, the OoD accuracy always has improved. Late stopping and tuning batch normalization momentum do not directly encourage to increase the SI score, but we observe that they do increase the SI score most of the cases, and when this happens, the OoD accuracy is also improved in more than 80.0 %of the cases. Thus, these results suggest that the improvements of OoD accuracy is strongly driven by the increase of the SI score.

Finally, we observe by analyzing $P(\Delta_{acc}^+|\Delta_{SI}^-)$, that when the SI score has not increased after applying one of the three approaches, the OoD accuracy still improves in a non-negligible number of cases. This suggests the existence of another mechanism that can improve the OoD accuracy even if the selectivity and invariance mechanisms did not emerge. However, one possible limitation of this interpretation is that selectivity and invariance may have emerged but have not been captured by the SI score, because the SI score may not quantify the emergence of these mechanisms in the most



Figure 6: *Correlation analysis*. This figure shows the correlation between OoD test accuracy and SI score. The Pearson's correlation coefficient is 0.891.

precise way. Thus, we can not make any assertion beyond the fact that it is unclear what are the neural mechanisms that facilitate OoD generalization when the three approaches do not manage to increase the SI score. This result motivates follow-up investigations.

In summary, in this study we provided evidence that the invariance and selectivity mechanism drives OoD generalization. Also, we found cases in which improvements of OoD generalization may not be preceded by the strengthening of the selectivity and invariance mechanism in the neural representations, which requires future work proposing novel mechanisms to explain this cases. We believe our experimental framework will facilitate these future discoveries.

7. Conclusion

We have shown that late-stopping, tuning the batch normalization momentum parameter, and optimizing the invariance loss during learning lead to substantial improvements of the DNN recognition accuracy of objects in OoD orientations and illuminations (in some cases more than 20%). These improvements are consistent across four datasets, and different degrees of dataset bias. We also corroborated that the neural mechanisms of selectivity to a category and invariance to orientations and illuminations, at the individual neuron level, lead to the aforementioned improvements of OoD recognition accuracy. Namely, we found that in the majority of trials where any of the three approaches yield an increase of selectiviy and invariance, resulted in improvements of the OoD recognition accuracy. Nonetheless, our analysis also revealed that other mechanisms different from selectivity and invariance may also exists, as we observed that gains of OoD recognition accuracy were not preceded by an increase of the SI score in

some trials. What are the neural mechanisms that drive OoD generalization in these cases remains as an open question for future work. Furthermore, there are also other novel questions derived from our results that motivate future works: Is there any effective way of combining the three approaches investigated in this paper that leads to even more improvements of OoD generalization? Are these approaches applicable to other factors beyond orientations and illumination conditions? How these approaches relate to biological learning systems? We hope that the substantial improvements of OoD recognition accuracy that we demonstrated in this paper motivate new research to address the fascinating questions that have cropped up ahead of us.

Finally, we would like to highlight that poor OoD generalization is one of the issues of machine learning that needs to be urgently addressed in order to allow for safe and fair AI applications. We hope that this research serves as a basis for further improvements of OoD generalization.

Acknowledgments

We are grateful to Tomaso Poggio and Hisanao Akima for their insightful advice and warm encouragement. We thank Shinichi Matsumoto and Shioe Kuramochi for their assistance to create CarsCG and DAISO-10 datasets, respectively. This work was supported by Fujitsu Limited (Contract No. 40008819 and 40009105) and by the Center for Brains, Minds and Machines (funded by NSF STC award CCF-1231216). PS and XB are supported by the R01EY020517 grant from the National Eye Institute (NIH).

Conflicts of Interests Statement

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

References

- A. Torralba, A. A. Efros, Unbiased look at dataset bias, in: CVPR, 2011, pp. 1521–1528.
- [2] R. Geirhos, P. Rubisch, C. Michaelis, M. Bethge, F. A. Wichmann, W. Brendel, ImageNet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness., in: ICLR, 2019.
- [3] S. Beery, G. V. Horn, P. Perona, Recognition in terra incognita, in: ECCV, 2018, pp. 472–489.

- [4] D. Hendrycks, S. Basart, N. Mu, S. Kadavath, F. Wang, E. Dorundo, R. Desai, T. Zhu, S. Parajuli, M. Guo, D. Song, J. Steinhardt, J. Gilmer, The many faces of robustness: A critical analysis of out-of-distribution generalization, arXiv preprint arXiv:2006.16241 (2020).
- [5] B. Recht, R. Roelofs, L. Schmidt, V. Shankar, Do CIFAR-10 classifiers generalize to CIFAR-10?, arXiv preprint arXiv:1806.00451 (2018).
- [6] B. Recht, R. Roelofs, L. Schmidt, V. Shankar, Do ImageNet classifiers generalize to ImageNet?, in: ICML, 2019, pp. 5389– 5400.
- [7] M. A. Alcorn, Q. Li, Z. Gong, C. Wang, L. Mai, W.-S. Ku, A. Nguyen, Strike (with) a pose: Neural networks are easily fooled by strange poses of familiar objects, in: CVPR, 2019, pp. 4845–4854.
- [8] A. Barbu, D. Mayo, J. Alverio, W. Luo, C. Wang, D. Gutfreund, J. Tenenbaum, B. Katz, ObjectNet: A large-scale biascontrolled dataset for pushing the limits of object recognition models, in: NeurIPS, 2019.
- [9] P. Sinha, T. Poggio, Role of learning in three-dimensional form perception, Nature 384 (6608) (1996) 460–463.
- [10] S. Ullman, High-level vision: Object recognition and visual cognition, MIT Press, 1996.
- [11] F. Anselmi, L. Rosasco, T. Poggio, On invariance and selectivity in representation learning, Information and Inference: A Journal of the IMA 5 (2) (2016) 134–158.
- [12] S. S. A. Zaidi, X. Boix, N. Prasad, S. Gilad-Gutnick, S. Ben-Ami, P. Sinha, Is robustness to transformations driven by invariant neural representations?, arXiv preprint arXiv:2007.00112 (2020).
- [13] S. Madan, T. Henry, H. Ho, N. Bhandari, T. Sasaki, F. Durand, H. Pfister, X. Boix, When and how do CNNs generalize to out-of-distribution category-viewpoint combinations?, arXiv preprint arXiv:2007.08032v2 (2021).
- [14] Y. Yao, L. Rosasco, A. Caponnetto, On early stopping in gradient descent learning, Constructive Approximation (2007) 289– 315.
- [15] S. Schneider, E. Rusak, L. Eck, O. Bringmann, W. Brendel, M. Bethge, Improving robustness against common corruptions by covariate shift adaptation, in: NeurIPS, 2020.
- [16] Y. LeCun, L. Bottou, Y. Bengio, P. Haffner, Gradient-based learning applied to document recognition, Proceedings of the IEEE 86 (11) (1998) 2278–2324.
- [17] A. Borji, S. Izadi, L. Itti, iLab-20M: A large-scale controlled object dataset to investigate deep learning, in: CVPR, 2016, pp. 2221–2230.
- [18] E. D. Cubuk, B. Zoph, J. Shlens, Q. V. Le, Randaugment: Practical data augmentation with no separate search, arXiv preprint arXiv:1909.13719 (2019).
- [19] S. S. Halder, J.-F. Lalonde, R. d. Charette, Physics-based rendering for improving robustness to rain, in: ICCV, 2019, pp. 10203–10212.
- [20] Y. Kim, A. F. M. S. Uddin, S. H. Bae, Local augment: Utilizing local bias property of convolutional neural networks for data augmentation, IEEE Access 9 (2021) 15191–15199.
- [21] F. Qiao, L. Zhao, X. Peng, Learning to learn single domain generalization, in: CVPR, 2020, pp. 12556–12565.
- [22] W. Angtian, A. Kortylewski, A. Yuille, NeMo: Neural mesh models of contrastive features for robust 3D pose estimation, in: ICLR, 2021.
- [23] B. Chidester, T. Zhou, M. N. Do, J. Ma, Rotation equivariant and invariant neural networks for microscopy image analysis, Bioinformatics 35 (14) (2019) i530–i537.
- [24] K. Qi, Q. Guan, C. Yang, F. Peng, S. Shen, H. Wu, Concentric circle pooling in deep convolutional networks for remote sens-

ing scene classification, Remote Sensing 10 (6) (2018).

- [25] S. Sabour, N. Frosst, G. E. Hinton, Dynamic routing between capsules, in: NeurIPS, 2017, pp. 3859–3869.
- [26] T. S. Cohen, M. Geiger, J. Köhler, M. Welling, Spherical CNNs, in: ICLR, 2018.
- [27] F. M. Carlucci, A. D'Innocente, S. Bucci, B. Caputo, T. Tommasi, Domain generalization by solving jigsaw puzzles, in: CVPR, 2019, pp. 2229–2238.
- [28] Q. Dou, D. C. Castro, K. Kamnitsas, B. Glocker, Domain generalization via model-agnostic learning of semantic features, in: NeurIPS, 2019, pp. 6450–6461.
- [29] M. Ghifary, W. Bastiaan Kleijn, M. Zhang, D. Balduzzi, Domain generalization for object recognition with multi-task autoencoders, in: ICCV, 2015, pp. 2551–2559.
- [30] J. Guo, X. Zhu, C. Zhao, D. Cao, Z. Lei, S. Z. Li, Learning meta face recognition in unseen domains, arXiv preprint arXiv:2003.07733 (2020).
- [31] Y. Jia, J. Zhang, S. Shan, X. Chen, Single-side domain generalization for face anti-spoofing, in: CVPR, 2020, pp. 8484–8493.
- [32] D. Li, Y. Yang, Y.-Z. Song, T. M. Hospedales, Learning to generalize: Meta-learning for domain generalization, in: AAAI, 2018.
- [33] H. Li, S. J. Pan, S. Wang, A. C. Kot, Domain generalization with adversarial feature learning, in: CVPR, 2018, pp. 5400–5409.
- [34] R. Volpi, H. Namkoong, O. Sener, J. C. Duchi, V. Murino, S. Savarese, Generalizing to unseen domains via adversarial data augmentation, in: NeurIPS, 2018.
- [35] P. Chattopadhyay, Y. Balaji, J. Hoffman, Learning to balance specificity and invariance for in and out of domain generalization, in: ECCV, 2020, pp. 301–318.
- [36] M. Ilse, J. M. Tomczak, C. Louizos, M. Welling, DIVA: Domain invariant variational autoencoder, in: MIDL, 2020, pp. 322–348.
- [37] R. Geirhos, J.-H. Jacobsen, C. Michaelis, R. Zemel, W. Brendel, M. Bethge, F. A. Wichmann, Shortcut learning in deep neural networks, Nature Machine Intelligence 2 (11) (2020) 665–673.
- [38] S. Sagawa*, P. W. Koh*, T. B. Hashimoto, P. Liang, Distributionally robust neural networks, in: ICLR, 2020.
- [39] M. Arjovsky, L. Bottou, I. Gulrajani, D. Lopez-Paz, Invariant risk minimization, arXiv preprint arXiv:1907.02893 (2019).
- [40] F. Ahmed, Y. Bengio, H. van Seijen, A. Courville, Systematic generalisation with group invariant predictions, in: ICLR, 2021.
- [41] E. Creager, J.-H. Jacobsen, R. Zemel, Environment inference for invariant learning, in: ICML, 2021, pp. 2189–2200.
- [42] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: CVPR, 2016, pp. 770–778.
- [43] Z. Cataltepe, Y. S. Abu-Mostafa, M. Magdon-Ismail, No free lunch for early stopping, Neural Computation 11 (4) (1999) 995–1009.
- [44] R. Caruana, S. Lawrence, L. Giles, Overfitting in neural nets: Backpropagation, conjugate gradient, and early stopping, NeurIPS (2001) 402–408.
- [45] Y. Yao, L. Rosasco, A. Caponnetto, On early stopping in gradient descent learning, Constructive Approximation 26 (2) (2007) 289–315.
- [46] L. Prechelt, Early stopping-but when?, in: Neural Networks: Tricks of the trade, 1998, pp. 55–69.
- [47] S. Ioffe, C. Szegedy, Batch normalization: Accelerating deep network training by reducing internal covariate shift, in: ICML, 2015, pp. 448–456.
- [48] P. Luo, X. Wang, W. Shao, Z. Peng, Towards understanding regularization in batch normalization, in: ICLR, 2019.
- [49] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, et al., Tensorflow: A system for large-scale machine learning, in: OSDI, 2016, pp. 265–283.

- [50] R. Hadsell, S. Chopra, Y. LeCun, Dimensionality reduction by learning an invariant mapping, in: CVPR, 2006, pp. 1735–1742.
- [51] G. E. Dahl, T. N. Sainath, G. E. Hinton, Improving deep neural networks for LVCSR using rectified linear units and dropout, in: ICASSP, 2013, pp. 8609–8613.
- [52] X. Glorot, Y. Bengio, Understanding the difficulty of training deep feedforward neural networks, in: AISTATS, 2010, pp. 249–256.
- [53] D. P. Kingma, J. Ba, Adam: A method for stochastic optimization, in: ICLR, 2015.

Appendix A. Details of datasets

Appendix A.1. MNIST-Positions

Starting with the MNIST dataset [16], which is available at http://yann.lecun.com/exdb/mnist/ (Last access: Oct. 1, 2020), we created a dataset of 42×42 pixels with nine numbers (0 to 8) by resizing images to 14×14 and placing them in one of 9 possible positions in a 3×3 empty grid. We call it MNIST-Positions. Fig. A.7 shows the all categories and positions of MNIST-Positions. In our experiments, the numbers are considered to be the object category set C and the positions where the numbers are placed is considered as the condition set N. Thus, it is written as #(C) = #(N) = 9.

Appendix A.2. iLab-Orientations

iLab-2M is a dataset created from iLab-20M dataset [17]: available at https://bmobear.github.io/ projects/viva/ (Last access: Oct. 20, 2020). The dataset consists of images of 15 categories of physical toy vehicles photographed in various orientations, elevations, lighting conditions, camera focus settings and backgrounds. It has 1.2M training images, 270K validation images, 270K test images, and each image is 256×256 pixels. We chose from the original iLab-2M dataset six categories — bus, car, helicopter, monster truck, plane, and tank as *C* and six orientations as \mathcal{N} . We call it iLab-Orientations. Fig. A.8 shows samples of the all categories and orientations of iLab-Orientations dataset.

Appendix A.3. CarsCG-Orientations

CarsCG-Orientations is a new dataset that consists of images of ten models of cars in various conditions rendered by Unreal Engine version 4.25.3; we plan to make this dataset publicly available. The conditions consist of ten orientations, three elevations, ten body colors, five locations and three time slots. Fig. A.9 shows the all car models (categories) and orientations (conditions) in the grid form. The details of these are as follows.

- Categories: CarsCG-Orientations dataset consists of images of the following cars Nissan XTrail[®], Volkswagen[®] Golf, BMW 2Series Coupe[®], Honda Odyssey[®], Toyota Prius[®], Mercedes Benz[®] A-Class, Lexus[®] LS, Mercedes[®] Benz E-Class, Toyota Yaris[®] and Volvo[®] V40 (See Fig. A.10). We used the whole car models as categories C. Therefore the number of categories is #(C) = 10 in the experiments conducted in this study.
- Orientations: During the rendering process, the virtual camera (camera actor) was rotated around the yaw axis of each car from 0 to 324 degrees in units of 36 degrees. Therefore, each car model appears in the images with ten different azimuth orientations. All orientations are shown in Fig. A.11. We used the whole orientations as conditions N. Thus the number of the conditions is #(N) = 10 in the experiments conducted in this study.

To create variety of samples for each combination of the categories (car models) and conditions (orientations), we added other conditions as follows.

- Elevations: The virtual camera was located at three elevation angles, namely, 10, 15, and 30 degrees, during the rendering process. Sample images taken from each angle are shown in Fig. A.12.
- Body colors: Each car model is rendered with ten colors, namely, black, light blue, green, red, white, beige, dark blue, orange, plum, and silver by using Automotive Materials (a library for Unreal Engine). Fig. A.13 shows sample images of Nissan XTrail rendered with these colors.
- Locations: We used a sample environment of an urban park contained in City Park Environment Collection. We chose five locations from the sample environment and modified them for our experiments. Sample images taken at each location are shown in Fig. A.14.
- Time slots: We used Ultra Dynamic Sky 3D model set to synthesize the three different times slots, namely, daytime, twilight, and midnight. Fig. A.15 shows the samples of these three time slots.

The number of images and the image size are as follows.

• Number of images and image size: The total number of images of this dataset is 45K = 10 (categories) $\times 10$ (orientations) $\times 3$ (elevations) $\times 10$ (body colors) $\times 5$ (locations) $\times 3$ (time slots). The images are rendered in 3840×2160 pixels and then resized to 1920×1080 pixels for the sake of anti-aliasing.

Appendix A.4. MiscGoods-Illuminations

MiscGoods-Illuminations is a subset of DAISO-10, a novel dataset constructed for this study; we plan to make this dataset publicly available. The dataset consists of images of ten physical miscellaneous goods taken with five illumination conditions, two object aspects, twenty object orientations, five camera angles. Figure A.16 shows the all miscellaneous goods (categories) and illumination conditions in the grid form. The details of these are as follows.

- Categories: As shown in Fig. A.16, DAISO-10 has ten types of miscellaneous goods stuffed dolphin, stuffed whale, metal basket, imitation plant, cup, cleaning brush, winding tape, lace yarn, bottled imitation tomatoes, and bottled imitation green apples. In this study, we selected the following five miscellaneous goods from DAISO-10 as the categories C stuffed dolphin, stuffed whale, metal basket, imitation plant and cup. Therefore the number of categories is #(C) = 5 in the experiments conducted in this study.
- Illumination conditions: As the conditions, we created five illumination conditions (lighting conditions); one is created with ceiling lights, and the rest are with a colored spotlight. All illumination conditions are shown in Fig. A.16. For spotlight conditions, the light source (PIXEL G1STM RGB Video Light) was placed 23 cm in front of the object (See Fig. A.17). The parameters of the light source were H217/S141=8500k (white light), H0/S100 (red light), H120/S100 (green light), and H240/S100 (blue light). These parameters were set so that the condition of the illumination makes a sufficient difference in the learning experiments. We used whole illumination conditions \mathcal{N} . Thus the number of the conditions is $\#(\mathcal{N}) = 5$ in the experiments conducted in this study.

As we did for CarCGs-Orientations, we added other conditions to create variety of samples for each combination of the categories and illumination conditions as follows.

- Object poses (aspects and orientations): In this dataset, we placed each object in two representative aspects for each lighting condition. Fig. A.18 shows the two aspects of all objects. For additional diversity, we rotated the object every 18 degrees from 0 to 342 degrees (Fig. A.19). In total, there are 40 patterns in object pose conditions.
- Camera angles: To capture the images automatically, we created a robotic image capture system (see Fig. A.17). A camera (Intel[®] Realsense D435) was attached to a robot arm (COBOTTA[®]), and the system captured images from five camera angles for each lighting and object pose condition (Fig. A.20). The postures were defined so that the acquired image shows the entire object pose. The series of operations from robot control to image acquisition is automated by utilizing ROS kinetic.

The number of images and the image size are as follows.

• Number of images and image size: The number of images of whole dataset is 10K = 10 (categories) \times 5 (illuminations) \times 2 (aspects) \times 20 (orientations) \times 5 (camera angles), and each image size is 640×480 pixels.

Appendix B. Details of experiments

ResNet-18 [42] is adopted as the network for all experiments. The source codes are implemented based on Python v3.6.9, using TensorFlow v2.5.0 and NumPy v1.19.5. They are included in the zip file (/source_code). The whole network architecture is shown in Fig B.1 and Fig B.2. All neurons employ the rectified linear function $g(z) = \max\{0, z\}$ and satisfy $a^{nm}(\mathbf{x}) \ge 0$. Glorot uniform initializer [52] is adopted for the network weights initialization for all experiments. We use BatchNormalization to standardize the inputs to a layer for each mini-batch. We use it for stabilizing the learning process and reducing the number of training epochs. We do not use any data augmentations. Invariance loss is applied to the last fully-connected layer "activation_17" with 512 neurons shown in Fig. B.1. Adam [53] is employed as the optimization algorithm. The cross-entropy loss is employed as the loss *L*. The pixels of images are

normalized within 0 to 1 as a preprocessing for all datasets. The epoch size and batch size are confirmed to produce reasonable accuracy in the baseline case for each dataset and we employ the same values for all experiments with the same dataset. For example, we use 100 and 256 as epoch size and batch size, respectively, for MNIST-Positions. The values of hyper-parameters are summarized in Table B.1. We have employed four Tesla V100 GPUs for the

Table B.1: Hyper-parameters used for each dataset

dataset	epoch size	preprocessing	weights initialization	batch size
MNIST-Positions	100	divide by 255	Glorot uniform initializer	256
iLab-Orientations	100	divide by 255	Glorot uniform initializer	256
CarsCG-Orientations	100	divide by 255	Glorot uniform initializer	32
MiscGoods-Illuminations	100	divide by 255	Glorot uniform initializer	32

experiments. The preparation of training dataset $\mathcal{D}_{train}^{(InD)}$, InD validation dataset $\mathcal{D}_{val}^{(InD)}$, and OoD dataset $\mathcal{D}^{(OoD)}$ has been conducted as follows.

- MNIST-Positions: We use images of the original MNIST-Positions with image size of 42×42 pixels. InD dataset and OoD dataset are prepared in the way described in Sec. 3.1. The number of train dataset is #(D^(InD)_{train}) = 54000. We use #(D^(InD)_{val}) = 8000 for InD validation dataset. The number of OoD dataset is #(D^(OoD)) = 8000.
- iLab-Orientations: We reize the images to 64×64 pixels. InD dataset and OoD dataset are prepared in the way described in Sec. 3.1.The number of train dataset is $\#(\mathcal{D}_{train}^{(InD)}) = 18000$. We use $\#(\mathcal{D}_{val}^{(InD)}) = 8000$ for InD validation dataset. The number of OoD dataset is $\#(\mathcal{D}^{(OoD)}) = 8000$.
- CarsCG-Orientations: We resize the images to 224×224 pixels. InD dataset and OoD dataset are prepared in the way described in Sec. 3.1.The number of train dataset is $\#(\mathcal{D}_{train}^{(InD)}) = 3400$. We use $\#(\mathcal{D}_{val}^{(InD)}) = 450$ for InD validation dataset. The number of OoD dataset is $\#(\mathcal{D}^{(OoD)}) = 800$.
- MiscGoods-Illuminations: We resize the images to 224×224 pixels. InD dataset and OoD dataset are prepared in the way described in Sec. 3.1. The number of train dataset is $\#(\mathcal{D}_{train}^{(InD)}) = 800$. We use $\#(\mathcal{D}_{val}^{(InD)}) = 200$ for InD validation dataset. The number of OoD dataset is $\#(\mathcal{D}^{(OoD)}) = 400$.

Appendix C. Additional results of experiments

InD accuracy and OoD accuracy learning curves with all dataset and all diversity corresponding to Fig. 1(a) are available in Fig. C.3 C.4 C.5. Furthermore InD accuracy and OoD accuracy learning curves with all all dataset and all diversity corresponding to Fig. 1(b) are available in Fig. C.6 C.7 C.8. InD accuracy corresponding to Fig. 5(a) 5(b) 5(c) is available in Fig. C.9. InD accuracy of difference from baseline corresponding to Fig. 5(d) 5(e) 5(f) is also available in Fig. C.10 The experiments for measuring accuracy is exactly same as what we reported in the main body of the paper.

Appendix D. Details of combined method

For the best of three approaches, we choose the one with the highest OoD accuracy calculated on the data set for hyperparameter tuning among late stopping, tuning batch normalization momentum, and invariance loss. The combined three approaches employs the hyper-parameters that determined in the Sec. 6: epoch size for longer epochs, momentum parameter β for tuning batch normalization momentum, and learning rate, paring interval and the value λ for invariance loss. If the InD accuracy drops to a chance, the learning rate is multiplied by 0.1 and the training is re-started.

0	0	0	0	0	0	0	0	0
/	/	/	/	/	/	/	/	/
2	2	2	a	2	a	2	a	2
3	3	3	3	3	3	3	3	3
Ч	ч	Ч	Ч	ч	ч	ч	ч	ч
5	ч	5	ч	Ч	5	5	5	5
6	6	6	6	6	6	6	6	6
7	7	7	7	7	7	7	7	7
8	8	8	8	8	8	8	8	8

Figure A.7: Sample images of MNIST-Positions dataset arranged in a grid pattern. Each row indicates the number as the object category. MNIST-Poitions include nine numbers from 0 to 8. Each column indicates the positions as the condition category. There are 9 positions in this dataset.



Figure A.8: Sample images of iLab–Orientations dataset arranged in a grid pattern. Each row indicates the object category. iLab-Orientations include six object categories — bus, car, helicopter, monster truck, plane, and tank. Each column indicates the orientations as the condition category. There are 6 orientations in this dataset.



Figure A.9: Sample images of each object category and orientation of CarsCG-Orientations. Each row indicates object categories — Nissan XTrail[®], Volkswagen [®] Golf, BMW[®] 2Series Coupe, Honda Odyssey[®], Toyota Prius[®], Mercedes Benz[®] A-Class, Lexus[®] LS, Mercedes Benz[®] E-Class, Toyota Yaris[®] and Volvo[®] V40. Each column indicates the condition categories, orientations from 0 to 324 degrees. These categories and orientations in this figure are used in our experiments.



Figure A.10: Sample images of ten object categories of CarsCG-Orientatoins — Nissan XTrail[®], Volkswagen[®] Golf, BMW[®] 2Series Coupe, Honda Odyssey[®], Toyota Prius[®], Mercedes Benz[®] A-Class, Lexus[®] LS, Mercedes Benz[®] E-Class, Toyota Yaris[®] and Volvo[®] V40 are shown in this figure from the top left to the bottom right. All conditions except object category are fixed in this figure.



Figure A.11: Sample images of ten orientations (condition categories) of CarsCG-Orientations. Ten orientations from 0 to 324 degrees are displayed from the left to right. All conditions except orientation are fixed in this figure.



Figure A.12: Sample images of the three elevation angles (10, 15, 30 degrees) of CarsCG-Orientations. Left figure is the image whose elevation angle is 10 degrees. Middle figure is the image whose elevation angle is 15 degrees. Right figure is the image whose elevation angle is 30 degrees.



Figure A.13: Sample images of cars painted in ten colors of CarsCG-Orientations. There are images painted in black, light blue, green, red, white, beige, dark blue, orange, plum, and silver from left to right.



(a) Sample images of each location from elevation angle of 10 degrees that CarsCG-Orientations has.



(b) Sample images of each location from elevation angle of 15 degrees that CarsCG-Orientations has.

Figure A.14: (a) is images of a car placed in five different locations taken from elevation angle of 10 degrees. (b) is images taken from elevation angle of 15 degrees. There are differences in texture of road and background.



Figure A.15: Sample images of each time slot of CarsCG-Orientations. Left figure is an image of car taken in the daytime. Middle figure shows an image of a car taken in the twilight. The color of the car is different from that in the left and right ones. It is caused by the twilight sunlight condition. Right figure shows an image of car taken at midnight. The color of the car is also different from left and middle ones.



Figure A.16: Sample images of each object category and illumination condition of MiscGoods-Illuminations are shown in this figure. Each row indicates object categories — stuffed dolphin, stuffed whale, metal basket, imitation plant, cup, cleaning brush, winding tape, lace yarn, bottled imitation tomatoes, and bottled imitation green apples. Each column indicates the condition categories, illumination conditions — ceiling light, white spotlight, red spotlight, green spotlight and blue spotlight. These five categories from the top and five illumination conditions are used as object categories and condition categories in our experiments.



Figure A.17: Robotic image capture system for MiscGoods-Illuminations. Dashed bidirectional arrow indicates the robot motion.



Figure A.18: Sample images of MiscGoods-Illuminations with two aspects. Each object has these two aspects as condition. The shapes of these objects in an image are changed by aspects conditions.



Figure A.19: Sample images of each orientation of MiscGoods-Illuminations. 20 orientations from 0 to 342 degrees that the dataset has are shown in this figure from the top left to the bottom right.



Figure A.20: Sample images from each camera angles of MiscGoods-Illuminations. There are five angles in the dataset. The postures were defined so that the acquired image shows the entire object pose. These five camera angles are related to postures of robot arm that the camera is connected.



Figure B.1: This diagram shows the whole architecture of our implementation of ResNet-18. The numbers in this diagram represent batchsize, height of image, width of image and channels. Therefore they change depending on the dataset. Current numbers correspond to MNIST-Positions. For instance, the numbers on the top of the diagram means (batchsize, height, width, channels) = (128, 42, 42, 1). Conv2D, Dense and BasicBlock mean a convolutional layer, a fully connected layer and a basic building block of ResNet, respectively.



Figure B.2: This diagram shows the architecture of BasicBlock in ResNet-18. Conv2D and Add mean a convolutional layser and a layer that simply add the two input values.



Figure C.3: Late-stopping with low InD data diversity



Figure C.4: Late-stopping with medium InD data diversity



Figure C.5: Late-stopping with high InD data diversity



Figure C.6: Baseline and BN momentum with low InD data diversity



Figure C.7: Baseline and BN momentum with medium InD data diversity



Figure C.8: Baseline and BN momentum with high InD data diversity



Figure C.9: Performance improvement in mean InD accuracy



Figure C.10: Performance improvement of difference from Baseline in mean InD accuracy



Figure C.11: Overall results between IS scores and all combinations of (data diversity, dataset, approaches)