

ROBUSTNESS MAY BE MORE BRITTLE THAN WE THINK UNDER DIFFERENT DEGREES OF DISTRIBUTION SHIFTS

Anonymous authors

Paper under double-blind review

ABSTRACT

Out-of-distribution (OOD) generalization is a complicated problem due to the idiosyncrasies of possible distribution shifts between training and test domains. Most benchmarks employ diverse datasets to address this issue; however, the degree of the distribution shift between the training domains and the test domains of each dataset remains largely fixed. This may lead to biased conclusions that either underestimate or overestimate the actual OOD performance of a model. Our study delves into a more nuanced evaluation setting that covers a broad range of shift degrees. We show that the robustness of models can be quite brittle and inconsistent under different degrees of distribution shifts, and therefore one should be more cautious when drawing conclusions from evaluations under a limited range of degrees. In addition, we observe that large-scale pre-trained models, such as CLIP, are sensitive to even minute distribution shifts of novel downstream tasks. This indicates that while pre-training may improve downstream in-distribution performance, it could have minimal or even adverse effects on generalization in certain OOD scenarios of the downstream task. In light of these findings, we encourage future research to conduct evaluations across a broader range of shift degrees whenever possible.

1 INTRODUCTION

Out-of-distribution (OOD) generalization is vital to the safety and reliability of machine learning applications in the real world. However, the complexities of distribution shifts between the training domains and the real test domains make OOD generalization a challenging problem. Numerous empirical studies (Gulrajani & Lopez-Paz, 2021; Wiles et al., 2022) have suggested that most algorithms only offer very little improvement in OOD performance over empirical risk minimization (ERM) (Vapnik, 1998). Furthermore, algorithms performing better than ERM against one type of distribution shift often perform poorly against another (Ye et al., 2022). The inconsistency suggests that it is important to consider various possible types of distribution shifts of a task when evaluating the OOD performance of a model; otherwise, the evaluation might lead to biased conclusions.

To address the issue, most OOD benchmarks (Koh et al., 2021; Hendrycks et al., 2021; Gulrajani & Lopez-Paz, 2021; Ye et al., 2022) incorporate multiple datasets exhibiting a diverse range of distribution shifts. However, another potential source of evaluation bias is often overlooked: the test domains only capture a largely fixed degree of each distribution shift. For example, in (Li et al., 2017; Koh et al., 2021; He et al., 2021; Zhao et al., 2022), each test domain represents a different “direction” of the potential distribution shifts of a task but there is no distinction between different degrees of shift on the same direction. Similar problems can also arise when only the aggregate performance across multiple degrees is examined (Hendrycks & Dietterich, 2019). Such kind of evaluation can result in misconceptions about model performance on the same grounds as those of the evaluation based on limited types (or “directions”) of distribution shifts.

Consider the situation (that we observed in this work) illustrated in Figure 1, where the performance of a model is evaluated in only two domains, one for in-distribution (ID) performance in the training domain $\mathcal{D}_{\text{train}}$, and the other for OOD performance in the test domain $\mathcal{D}_{\text{test}}$. In this case, the observed

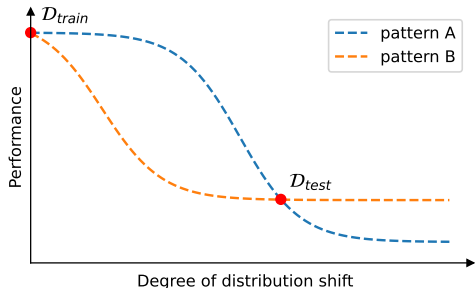


Figure 1: A typical situation where an evaluation under a limited number of degrees of a distribution shift cannot tell any difference between two distinct OOD generalization patterns (labeled as A and B) that can be realized by a model.

performance, which can be explained by at least two distinct generalization patterns as shown in the figure, presents an oversimplified summary of the OOD generalization ability of the model. This simplification may lead to incorrect assumptions about model robustness under various degrees. For example, when a model outperforms another model under a distribution shift of certain severity, it could leave the wrong impression that the first model is more robust in general, i.e., outperforming the other model under almost every possible degree of the concerned distribution shift, while in fact, the first model has much poorer worst-case performance.

In this study, we take a closer look at OOD generalization under distribution shifts of varying degrees. We are interested in the behavior of different models under a broad range of shift degrees and also the relation between the performance of a model at different degrees. Through extensive experiments, we make several observations about the generalization behavior of models under the considered evaluation setting. First, we highlight that the advantage of a model under a mild shift may not apply to stronger shifts of the same type, even if the shift is just slightly stronger¹. Therefore, caution should be taken when interpreting evaluation results obtained under a limited range of shift degrees. Second, we find that training a model with strongly shifted data can sometimes guarantee robustness to all milder shifts, while at other times it only has a limited impact on robustness and may even harm the OOD performance under milder shifts.

Lastly, the brittleness of robustness to different degrees of distribution shift is also observed in large-scale pre-trained models. We find that while CLIP (Radford et al., 2021) models are able to adapt to many novel tasks, achieving great (sometimes near-perfect) downstream ID performance, they can be extremely sensitive to downstream distribution shifts. In the presence of a distribution shift that is rarely seen during pre-training, even a very mild degree of the shift can cause a disproportionate performance drop in CLIP models in comparison to models trained from scratch. Interestingly, further adapting to the shift to which the models are sensitive significantly improves their general robustness. We believe that such characterizations of the “growth” of the OOD generalization ability of a model is a generally good practice. We encourage future research to adopt this kind of evaluation to generate more valuable insights into OOD generalization.

2 RELATED WORK

Out-of-distribution (OOD) generalization. Deep neural networks have demonstrated incredible generalization on a variety of complicated tasks, sometimes exceeding human performance (Rusakovsky et al., 2015; Silver et al., 2017; OpenAI, 2023), but they are shown to generalize very differently as we do and are very sensitive to all kinds of distribution shifts (Szegedy et al., 2014; Geirhos et al., 2020; Wang et al., 2023). Such brittleness severely undermines the reliability of neural networks and hence limits their applications in the real world where the stakes can be very high. For this reason, OOD generalization and related areas such as domain generalization (Blan-

¹We use mild/strong and low/high-degree interchangeably when describing a distribution shift.

chard et al., 2011; Zhou et al., 2021; Wang et al., 2022) has gained much attention rapidly in recent years (Shen et al., 2021).

Distribution shifts and OOD benchmarks. A considerable amount of the research efforts on OOD generalization has been dedicated to the evaluation of models and methods. Hendrycks & Dietterich (2019) provided benchmarks for evaluating the robustness of deep models against common image corruptions and perturbations. One of the benchmarks, ImageNet-C, consists of images under different severity levels of corruptions. In this benchmark, the authors examined the average accuracy of a number of models over all severity levels and showed that the models were all vulnerable to the considered corruptions. They did not, however, provide any analysis at the level of each individual severity level of corruption. Hendrycks et al. (2021) further proposed OOD benchmarks under natural distribution shifts, but this time like DG benchmarks such as (Gulrajani & Lopez-Paz, 2021), they do not involve any evaluation or discussion with regard to different degrees of distribution shifts. Meanwhile, Koh et al. (2021) proposed a diverse set of OOD benchmarks derived from real-world tasks but the degrees of distribution shifts are still largely fixed in each task. Similar examples include (Peng et al., 2019; He et al., 2021; Liang & Zou, 2022; Zhao et al., 2022) which focus on incorporating as many diverse types of distribution shifts without considering different degrees of distribution shifts. Lynch et al. (2023) considered three levels of spurious correlation but did not discuss the connection between the model performance at each level.

Distribution shifts in model learning. Schott et al. (2022) showed that models regardless of supervision signal and architectural bias could not learn the underlying mechanism that causes the distribution shifts on several datasets with controllable factors. In comparison, our finding suggests that learning the shifting-inducing mechanism of certain task is possible if the model can be made to be robust to the highest possible degree of the distribution shift. In a different context, Shi et al. (2022) also studied OOD generalization under multiple degrees of distribution shift. Their main focus is whether unsupervised methods can learn more robust representations than supervised learning. They conducted evaluation against three different degrees of spurious correlation and found that unsupervised methods are generally more robust than supervised learning and the advantage grows as the degree of the distribution shift increases.

Robustness of foundation models. Pre-training usually has a great impact on generalization. Foundation models such as CLIP (Radford et al., 2021) leverage massive scale of training data to generalize to a great variety of downstream tasks. At the same level of ID accuracy, zero-shot CLIP models are able to attain much higher OOD accuracy on several ImageNet variants than other models trained with a much smaller scale of data. Despite the breakthrough, however, the authors of CLIP have cautioned that the zero-shot performance of CLIP is significantly worse than existing models on data that are hardly present in the training data, e.g., MNIST (LeCun et al., 2010). Later, it is further shown that the main source of the remarkable robustness of CLIP is the diversity of its training data distribution (Fang et al., 2022). While CLIP can be made even more robust in some tasks after proper adaptation (Wortsman et al., 2022), what remains unclear in the literature is to what extent the robustness of CLIP and other foundation models can transfer to downstream tasks and how the models would behave as the degree of the downstream distribution shifts increases.

3 DIFFERENT DEGREES OF DISTRIBUTION SHIFTS

The degree of a distribution shift can be quantified in many ways. In this paper, we do not restrict our study to a particular way of quantification as different types of distribution shift may favor different ways of quantification. Instead, we consider the degrees of only one type of distribution shift at a time, so the degrees can take arbitrary values as long as they preserve a certain ordering of a set of domains under the same type of distribution shift.

For simplicity, we use natural numbers to represent the order of a given set of domains under the same type of distribution shift, with smaller numbers indicating lower degrees of distribution shift. We use \mathcal{D}_d to denote a domain where d is its degree and refer to \mathcal{D}_0 as a clean domain or a domain under no distribution shift. A model is said to be more robust to a degree d of distribution shift than another model if it attains better performance in \mathcal{D}_d .



Figure 2: Examples of the NOISYMNIST dataset which consists of 11 subsets of MNIST, one of which is clean, while the other 10 subsets are affected by different degrees of Gaussian noise.

A simple example of problems under different degrees of distribution shift is shown in Figure 2. In this example, the degree of distribution shift in a domain corresponds to the intensity of pixel-level Gaussian noise in the images. In practice, we may not have access to data under all possible degrees of a distribution shift and must rely on data under a limited set of shift degrees to train and evaluate our models. For this context, out-of-distribution generalization refers to the generalization from data under a given set of shift degrees to the rest of possible degrees.

This setting allows for a nuanced understanding of how different degrees of distribution shifts impact the learning and generalization capabilities of models. Although we focus on image classification problems and use accuracy as the metric for measuring model performance, we believe that the conclusions drawn from our experiment results also hold for more general problems.

4 ROBUSTNESS MAY BE MORE BRITTLE THAN WE THINK

This section explores the complexities and nuances of model robustness under varying degrees of distribution shifts in neural networks. We first demonstrate that models exhibiting robustness under a certain degree of shift can experience substantial performance degradation under slightly higher degrees of shift. Then we explore whether models robust to high degrees of shift maintain this robustness under lower ones, revealing contrasting results dependent on the specific task and dataset, thereby highlighting the inherent brittleness in neural network robustness under different degrees of distribution shift.

4.1 EXPERIMENT SETUP

Datasets. Our study employs two altered versions of the MNIST dataset (LeCun et al., 2010), herein referred to as NOISYMNIST and ROTATEDMNIST. To introduce varying degrees of distribution shifts, the NOISYMNIST dataset is generated by introducing Gaussian noise to the original images, resulting in 10 shifted domains under different degrees. More specifically, The standard deviation of the noise is linearly spaced between 0 and 0.8, in increments of 0.08, at the pixel level, normalized to the pixel value range of 0 to 1. Any pixel value beyond this range is clipped to fit within the 0-1 boundary. The ROTATEDMNIST dataset is created by rotating the original images, with degrees linearly spaced from 0 to 80, at intervals of 10 degrees, resulting in 8 shifted domains. Note that our ROTATEDMNIST is different from the ones in other papers, e.g., (Gulrajani & Lopez-Paz, 2021) which covers a smaller set of rotation degrees.

Our study also employs an altered version of CIFAR10 (Krizhevsky et al., 2009). The dataset is called LOWLIGHTCIFAR10. The distribution shift in this dataset is a combination of two primitive types of distribution shifts which are shifts in brightness and shot-noise intensity. Since photos captured in darker environments tend to exhibit more intense shot noises, this dataset simulates realistic photographic effects in photos captured under low-light conditions and hence is much more realistic than the MNIST variants.

Algorithms. We experiment with Empirical Risk Minimization (ERM, Vapnik (1998)) and over 20 Domain Generalization (DG) algorithms, including but not limited to Invariant Risk Minimization (IRM, Arjovsky et al. (2019)), Variance Risk Extrapolation (VREx, Krueger et al. (2020)), Spectral Decoupling (SD, Pezeshki et al. (2021)), Deep Correlation Alignment (CORAL, Sun et al. (2017)), Group Distributionally Robust Optimization (GroupDRO, Sagawa et al. (2020)), Representation Self Challenging (RSC, Huang et al. (2020)), Domain-Adversarial Neural Networks (DANN, Ganin et al. (2016)), Inter-domain Mixup (Mixup, Yan et al. (2020)), Adaptive Risk Minimization (ARM, Zhang et al. (2020)), and many others representing various OOD research areas (see the full list in Appendix A.2).

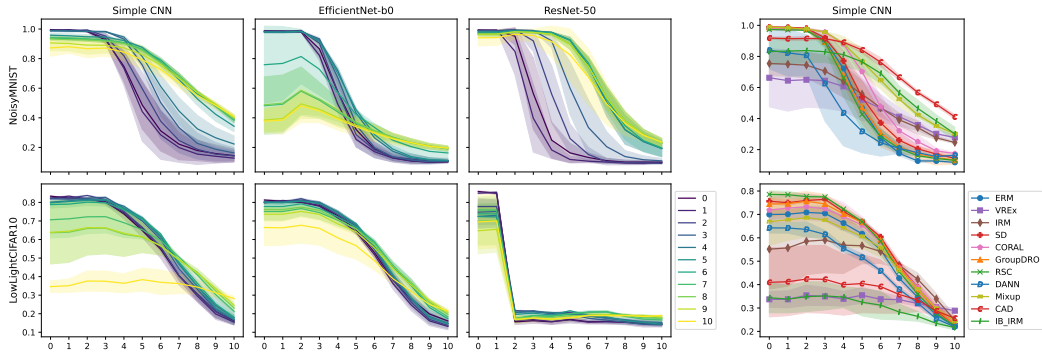


Figure 3: **(Left)** Performance of the best-performing models at each degree of NOISYMNIST and LOWLIGHTCIFAR10. The label of the curves denotes the domain on which the models perform best. The results are averaged over the top 5 models of all algorithms at each degree. **(Right)** Performance of ERM and representative domain generalization algorithms on the two datasets. The results are averaged over the top 3 models of each algorithm selected by worst-domain accuracy.

Implementation details. To conduct our experiments, we employed two neural network architectures: a simple 4-layer Convolutional Neural Network (CNN) and a more complex ResNet-50 (He et al., 2016) model. Both models were implemented without any form of pretraining. For optimization purposes, we utilized the Adam optimizer with a static learning rate of 0.001. The total batch size was fixed at 64, and was evenly divided across each training domain, and no weight decay was applied during the training process. Training iterations were set to a maximum of 5,000 for the 4-layer CNN and 10,000 for the ResNet-50 to ensure convergence. No form of data augmentation was used throughout the training process, preserving the inherent distribution and characteristics of the datasets. To ensure the reliability of our results, we conducted a thorough random search for hyperparameters, repeated 20 times.

4.2 ROBUSTNESS MAY NOT EVEN EXTRAPOLATE TO SLIGHTLY HIGHER DEGREES

In the real world, data under severe distribution shifts are usually very rare. We often face situations where we only have access to a reasonable amount of data under relatively mild distribution shifts. With these data, we can be fairly certain about the performance of a model in mild situations, but this is hardly satisfactory for any application that demands a certain level of reliability also in worse situations. Therefore, an important question is: how much can the performance of a model under some distribution shift tell us about its performance under stronger shifts?

To approach the question, we constructed a dataset, NOISYMNIST, by gradually adding Gaussian noise to MNIST (LeCun et al., 2010). As illustrated in Figure 2, NOISYMNIST consists of a clean subset \mathcal{D}_0 of MNIST and 10 subsets $\{\mathcal{D}_i\}_{i=1}^{10}$ under different degrees of noise. While the construction process of NOISYMNIST is simple, the dataset is nonetheless representative of a wide range of distribution shifts that gradually corrupt predictive features in an image. Intuitively, models that are more robust to relatively mild noises should also be more robust to stronger noises, at least to some extent. If this is true, then in the case of NOISYMNIST, we should be able to rely on domains under only mild shifts such as \mathcal{D}_4 , to pick the best-performing models in slightly worse domains such as \mathcal{D}_5 and \mathcal{D}_6 or even much worse domains such as \mathcal{D}_{10} .

For this investigation, we trained a pool of models on \mathcal{D}_0 and \mathcal{D}_1 with ERM and more than 20 domain generalization (DG) algorithms. The models share the same architecture (a 4-layer CNN) but are trained with different initializations and hyperparameters in addition to the different learning algorithms. The performance of the best-performing models in each domain is shown in Figure 3 (left, ERM+DG). The result indicates that models that are better under milder shifts are often significantly *worse* than the other models under stronger shifts. In particular, the average accuracy of the best-performing models in \mathcal{D}_4 has dropped by more than 10% in \mathcal{D}_5 which is only under a slightly more intense noise than \mathcal{D}_4 .

Algorithm	CNN			ResNet-50		
	\mathcal{D}_4	\mathcal{D}_5	\mathcal{D}_6	\mathcal{D}_4	\mathcal{D}_5	\mathcal{D}_6
ERM	77.8±2.8 (0.0)	47.7±5.2 (38.7)	26.5±5.0 (66.0)	97.4±0.3 (0.0)	84.0±5.5 (13.8)	54.8±14.6 (43.8)
VREx	90.1±1.7 (0.0)	74.3±5.6 (17.6)	53.4±6.4 (40.8)	64.6±1.7 (0.0)	32.1±2.3 (50.3)	17.4±0.9 (73.1)
IRM	78.7±2.4 (0.0)	57.6±8.2 (26.8)	38.0±11.3 (51.7)	95.7±0.8 (0.0)	82.4±1.3 (13.9)	56.6±6.8 (40.9)
SD	81.7±1.2 (0.0)	57.7±2.3 (29.4)	35.7±2.1 (56.4)	97.8±0.4 (0.0)	92.4±2.1 (5.5)	76.5±6.0 (21.7)
GroupDRO	74.0±1.7 (0.0)	50.3±5.7 (32.1)	29.9±8.4 (59.6)	82.4±9.0 (0.0)	53.1±17.7 (35.5)	30.5±10.4 (63.0)
RSC	84.3±4.6 (0.0)	61.4±7.2 (27.2)	39.6±7.1 (53.0)	88.4±4.0 (0.0)	64.6±8.6 (26.9)	39.2±8.0 (55.6)
Mixup	93.2±0.4 (0.0)	84.1±2.3 (9.7)	69.2±1.9 (25.7)	85.4±3.7 (0.0)	49.2±16.2 (42.4)	26.6±13.2 (68.8)
CAD	94.1±1.0 (0.0)	78.7±3.1 (16.3)	58.6±4.0 (37.7)	78.8±19.6 (0.0)	50.6±22.0 (35.8)	30.5±13.1 (61.4)
IB-IRM	86.1±5.6 (0.0)	68.9±9.7 (19.9)	54.6±13.3 (36.5)	91.0±2.9 (0.0)	59.5±12.0 (34.6)	30.1±12.3 (66.9)

Table 1: Performance of the best models in \mathcal{D}_4 of ERM and representative DG algorithms. The relative performance drops (%) with respect to the performance in \mathcal{D}_4 are shown in the parentheses. All results are averaged over the top 3 models among 20 models with different initialization and hyperparameters for training.

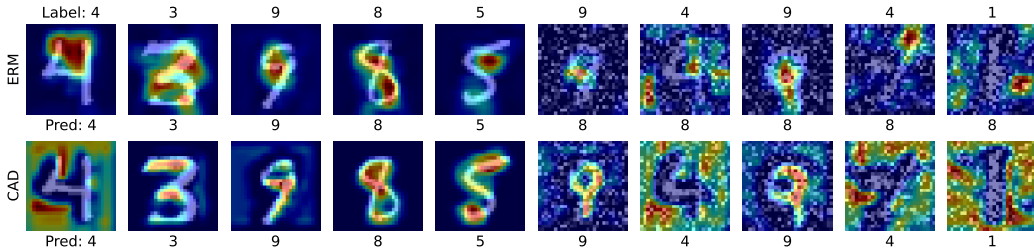


Figure 4: GradCAM visualization of model attention on random examples from \mathcal{D}_0 (left) and \mathcal{D}_7 (right) of NOISYMNIST. The two models (ERM and CAD) demonstrate distinctive generalization patterns, one relying on the local features while the other more on the global structures. The local features become unreliable as the noise becomes intense.

We further experimented with ResNet-50 to see if networks with much greater capacities can learn a representation that is generally more robust under all the considered levels of noise. As shown in Figure 3 (left, ERM+DG (ResNet-50)), the overall pattern still remains, although the difference among the best-performing models under the stronger end of noises has become less significant. While this shows that larger networks helps, the gap may never be able to be closed by increasing the capacity of the network. More importantly, *the robustness of a model may be more brittle than we think: even under the same type of distribution shift, a slight increase in the degree of the shift may severely harm the performance of the model.*

Besides individual models, the brittleness of robustness under different shift degrees also has implications in evaluating different learning algorithms. The performance of ERM and representative DG algorithms on NOISYMNIST are shown in Figure 3 (right), where the algorithms exhibit very different generalization patterns that cannot be accurately captured by evaluations under only a limited set of shift degrees. A number of DG algorithms, like ERM, are highly robust to low degrees of distribution shift but are quite brittle in the presence of higher degrees of shift. In contrast, there are also algorithms that are significantly better than ERM under high shift degrees but are worse than ERM in other cases. Moreover, when looking at the best-performing models in \mathcal{D}_4 , the same brittleness can be generally observed for all the algorithms in Table 1, where the performance drop can be even more drastic than that is shown in Figure 3. Astoundingly, the relative performance drop can go up to 50.3% from \mathcal{D}_4 to \mathcal{D}_5 and 73.1% from \mathcal{D}_4 to \mathcal{D}_6 .

To better understand the observed brittleness, we visualized the attention of ERM and CAD models on NOISYMNIST using GradCAM (Selvaraju et al., 2017). As shown in Figure 4, while both ERM and CAD models can make accurate predictions in the clean domain \mathcal{D}_0 , they rely on radically different patterns to do so. ERM prefers the most predictive features regardless of whether they are robust or not. In the case of NOISYMNIST, these features turn out to be local features, which are easily corrupted by the noise, and thus no longer predictive when the noise becomes intense.

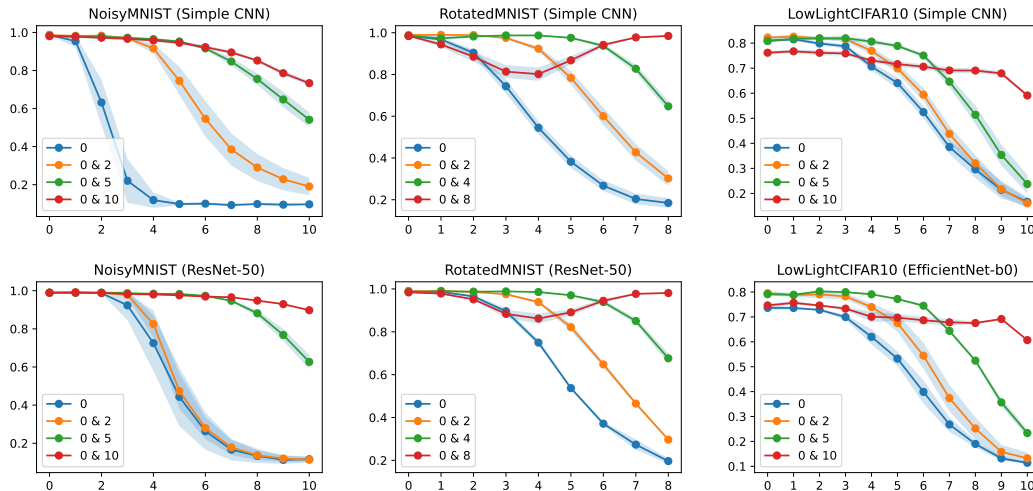


Figure 5: Performance of ERM models trained on domains under different shift degrees. The label of the curves denotes the indices of the training domains, e.g., “0 & 2” means that the models are trained on \mathcal{D}_0 and \mathcal{D}_2 of the corresponding dataset. The results are averaged over 20 models with different initialization.

From this perspective, we can see that the brittleness manifests when the spurious correlation between the local features and the target labels reaches a breaking point. However, where this breaking point is and how rapidly the correlation breaks seem to be totally dependent on the nature of the distribution shift and the task itself. While NOISYMNIST demonstrates a simple case where the breaking point is at a moderate degree of distribution shift, there can be scenarios where the break happens at a much lower or higher degree of distribution shift and happens much more rapidly. As a consequence, evaluations that only consider a narrow range of possible shift degrees would be highly unreliable in those scenarios.

4.3 ROBUSTNESS AT HIGHER DEGREES DOES NOT ALWAYS GUARANTEE ROBUSTNESS AT LOWER DEGREES

We have demonstrated in the previous section that models being more robust to milder distribution shifts does not imply that they are also more robust to stronger distribution shifts, even when the shift is just slightly stronger. In this section, we shed some light on the reverse question: does models being more robust to stronger distribution shifts imply them being more robust to milder distribution shifts?

To start with, we obtained models that are robust to strong distribution shifts by training the models on strongly shifted data together with clean data. In Figure 5, we compare these models with (i) models trained on much more mildly shifted data (also in addition to clean data) and (ii) models trained on clean data alone. For NOISYMNIST, the answer to our question is affirmative. The models that are more robust to stronger shifts are indeed more robust to milder shifts in general. However, this is not the whole story as the pattern seems to depend on the specific task in consideration.

We experimented with another dataset, ROTATEDMNIST, which is constructed in a similar fashion to NOISYMNIST while replacing noise with rotation (see Figure 7 for some examples). On the contrary to NOISYMNIST, robustness against higher shift degrees does *not* guarantee robustness to lower degrees in the case of ROTATEDMNIST when the shift degree of the additional training data is high. In comparison with models that trained on only clean data without any rotation, being more robust to the strongest shift may even harm generalization at the mildest degrees. Notably, the results are largely consistent between the two model architectures which have a great difference in complexity; however, the improvement brought by model complexity is still far from closing the gap. Again, this demonstrates the brittleness of the robustness of neural networks.

An important practical implication of the above finding is that, even for the same type of distribution shift, training on strongly shifted data may not be sufficient to obtain a model that is robust to milder shifts. Combined with our finding in Section 4, we arrive at the conclusion that the corresponding training data may be necessary to guarantee robustness at a certain degree of distribution shift for some tasks. Meanwhile, we should also note that there are scenarios where obtaining a dataset under a sufficiently strong shift is able to guarantee robustness to all milder shifts as in the case of NOISYMNIST. For these kinds of distribution shifts, it may require much less data to guarantee general robustness.

5 PRE-TRAINED REPRESENTATIONS ARE SENSITIVE TO NOVEL DOWNSTREAM DISTRIBUTION SHIFTS

Pre-training on large-scale datasets is one of the most effective ways that are known to consistently improve the generalization of neural networks across a wide range of tasks (Taori et al., 2020; Miller et al., 2021). In particular, foundation models like CLIP (Radford et al., 2021) have demonstrated remarkable zero-shot capability on a number of datasets that models trained on much smaller datasets fail to generalize to. In this section, we further investigate how pre-training would influence the OOD generalization behavior of models on downstream tasks under multiple shift degrees.

5.1 EXPERIMENT SETUP

Datasets. In addition to NOISYMNIST and ROTATEDMNIST, we consider two more complicated datasets, NOISYIMAGENET15 and LR-IMAGENET15, which are modifications of a 15-category subset of ImageNet on bird species. NOISYIMAGENET15 follows a similar construction to NOISYMNIST, introducing Gaussian noise on the pixel level, linearly spaced between 0 and 0.8, with values clipped to the 0-1 range. Meanwhile, LR-IMAGENET15 involves altering image resolution, first downsampling via bilinear interpolation and subsequently upsampling to 256×256 , with the downsampled resolution in each domain corresponding to a factor of $0.8^d \cdot 256$, where d represents the degree of distribution shift. We extend ROTATEDMNIST to span from 0 to 100 degrees in the experiments of this section.

Pre-trained models. ImageNet pre-trained ResNet-50 and ViT-B/32 from torchvision, along with CLIP checkpoints of these models released by OpenAI, serve as our primary models. These are adapted to downstream tasks through linear probing aligned with (Radford et al., 2021).

Implementation details. We use training-domain validation to select the best models among different iterations. All MNIST-based datasets were resized to 224×224 for uniformity across models. Pre-trained models normalized all datasets based on the statistics of their respective pre-training datasets, while randomly initialized models normalized based on MNIST statistics. Specifically, for the randomly initialized ResNet-50, no data augmentation was implemented during training on MNIST-based datasets. For the ViT-B/32 model, random affine transformations were applied, with rotation and shearing disabled for ROTATEDMNIST. We do not use any data augmentation for the experiments on NOISYIMAGENET15 and LR-IMAGENET15.

5.2 RESULTS

If large-scale pre-trained models like CLIP have learned generally more robust visual representations, they should be able to support a classifier that has better performance across a broad range of shift degrees than models trained from scratch as well as models trained on much smaller datasets, e.g., ImageNet (Deng et al., 2009). In Figure 6, we compare CLIP models with ImageNet pre-trained models and randomly initialized models that are trained from scratch on the downstream tasks to investigate the robustness of pre-trained representations.

On NOISYMNIST, although the pre-trained models perform equally well on clean domains as the randomly initialized models, they are surprisingly much more brittle to the distribution shift induced by the noise. Notably, the gap of accuracy between $CLIP_0$ and RI_0 increased by more than 40% from \mathcal{D}_0 to \mathcal{D}_1 on ResNet-50. This gap continued to increase until the shift reached a moderate degree. Moreover, on ViT-B/32, a similar pattern is observed albeit slightly improved. We hypothesize

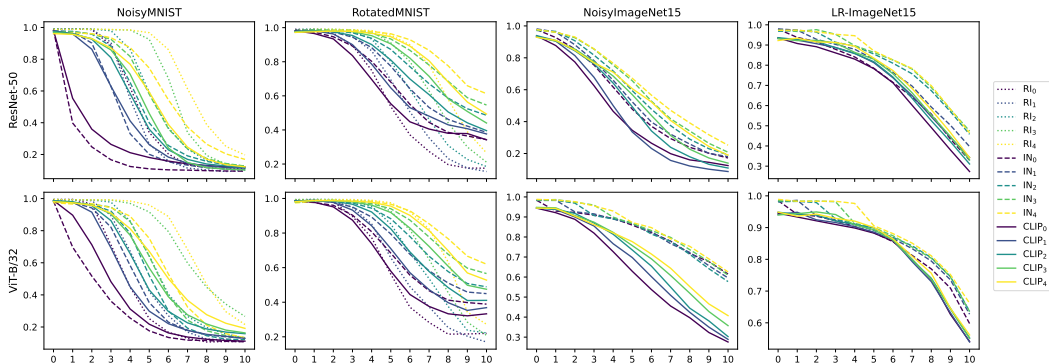


Figure 6: Performance of randomly initialized (**RI**) models, ImageNet (**IN**) pre-trained models, and **CLIP** models on different downstream tasks, evaluated over a broad range of shift degrees. The color of the curves indicates the domains used to train/adapt the models, e.g., RI_d stands for models trained on $\{\mathcal{D}_0, \dots, \mathcal{D}_d\}$ from scratch. The pre-trained models are adapted to the downstream tasks through linear probing. The results are averaged over three runs. Error bars are omitted for clarity (see Appendix B.5 for more details).

that the sensitiveness is largely because Gaussian noise is very rare in the training data of CLIP and also in ImageNet. Evaluation under a more common type of distribution shift, rotation, has provided some evidence to support our hypothesis. On **ROTATEDMNIST**, the pre-trained models are only slightly worse than the randomly initialized models under mild to moderate shifts while being much better under strong shifts.

In Figure 6, we also compare CLIP pre-trained models with ImageNet pre-trained models on harder problems: **NOISYIMAGENET15** and **LR-IMAGENET15**. On these two datasets, we observe that ImageNet pre-trained models are generally more robust than CLIP models under both distribution shifts. Furthermore, the gap between the two model families starts out being small but gradually enlarges as the shift gets stronger. This suggests that not only the property of the downstream distribution shift (e.g., noise) but also the difference between the pre-training task and the downstream task itself plays a role in determining the robustness of the pre-trained models against downstream distribution shifts. This implies that even if the pre-training data encompass a diverse set of distribution shifts, the robustness against those shifts may not transfer or only transfer very little to a different task under similar kinds of shifts.

Last but not least, we note that further adapting the pre-trained models to downstream distribution shifts can sometimes significantly improve their robustness as shown in Figure 6. On one hand, this corroborates existing findings that large-scale pre-trained representations are highly versatile. On the other hand, this also suggests that unleashing the power of pre-training may still require sufficiently diverse downstream task data that covers the potential distribution shifts. Nevertheless, there are still inherent limits to the power of pre-training under novelty downstream distribution shifts as demonstrated in the case of **NOISYMNIST** where further adaptations help but only to a limited degree compared with training from scratch.

6 CONCLUSION

In this work, we have shown that even when a model is robust to some degree of distribution shift, a slight increase in the degree of the shift can still cause a significant performance drop. In addition, we also find that training with data under a high degree of distribution shift sometimes guarantees robustness to all lower degrees, but not always. Furthermore, we observe that large-scale pre-trained models like CLIP are sensitive to downstream distribution shifts, especially unseen or rarely seen ones. These findings all suggest that the robustness of neural networks under certain degrees of distribution shift can be quite brittle. For this reason, we should be more careful when interpreting evaluations based on data with a limited range of shift degrees. We also encourage future research to adopt a more comprehensive evaluation of OOD generalization such as the one in this work that considers multiple degrees of distribution shifts.

REFERENCES

- Kartik Ahuja, Ethan Caballero, Dinghuai Zhang, Jean-Christophe Gagnon-Audet, Yoshua Bengio, Ioannis Mitliagkas, and Irina Rish. Invariance principle meets information bottleneck for out-of-distribution generalization. *Advances in Neural Information Processing Systems*, 34:3438–3450, 2021.
- Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization. *arXiv:1907.02893*, 2019.
- Gilles Blanchard, Gyemin Lee, and Clayton Scott. Generalizing from several related classification tasks to a new unlabeled sample. In *NeurIPS*, 2011.
- Mathieu Chevalley, Charlotte Bunne, Andreas Krause, and Stefan Bauer. Invariant causal mechanisms through distribution matching. *arXiv preprint arXiv:2206.11646*, 2022.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Cian Eastwood, Alexander Robey, Shashank Singh, Julius Von Kügelgen, Hamed Hassani, George J Pappas, and Bernhard Schölkopf. Probable domain generalization via quantile risk minimization. *Advances in Neural Information Processing Systems*, 35:17340–17358, 2022.
- Alex Fang, Gabriel Ilharco, Mitchell Wortsman, Yuhao Wan, Vaishal Shankar, Achal Dave, and Ludwig Schmidt. Data determines distributional robustness in contrastive language image pre-training (clip). In *International Conference on Machine Learning*, pp. 6216–6234. PMLR, 2022.
- Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *JMLR*, 2016.
- Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A. Wichmann. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2020.
- Ishaan Gulrajani and David Lopez-Paz. In search of lost domain generalization. In *ICLR*, 2021.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Yue He, Zheyang Shen, and Peng Cui. Towards non-iid image classification: A dataset and baselines. *Pattern Recognition*, 110:107383, 2021.
- Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. In *ICLR*, 2019.
- Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical analysis of out-of-distribution generalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 8340–8349, 2021.
- Zeyi Huang, Haohan Wang, Eric P Xing, and Dong Huang. Self-challenging improves cross-domain generalization. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*, pp. 124–140. Springer, 2020.
- Daehee Kim, Youngjun Yoo, Seunghyun Park, Jinkyu Kim, and Jaekoo Lee. Selfreg: Self-supervised contrastive regularization for domain generalization. In *ICCV*, 2021.
- Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Bal-subramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, et al. Wilds: A benchmark of in-the-wild distribution shifts. In *International Conference on Machine Learning*, pp. 5637–5664. PMLR, 2021.

- Masanori Koyama and Shoichiro Yamaguchi. Out-of-distribution generalization with maximal invariant predictor. *arXiv:2008.01883*, 2020.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- David Krueger, Ethan Caballero, Joern-Henrik Jacobsen, Amy Zhang, Jonathan Binas, Remi Le Priol, and Aaron Courville. Out-of-distribution generalization via risk extrapolation (rex). *arXiv:2003.00688*, 2020.
- Yann LeCun, Corinna Cortes, and CJ Burges. Mnist handwritten digit database. *ATT Labs [Online]*. Available: <http://yann.lecun.com/exdb/mnist>, 2, 2010.
- D. Li, Y. Yang, Y. Song, and T. M. Hospedales. Deeper, broader and artier domain generalization. In *ICCV*, 2017.
- Haoliang Li, Sinno Jialin Pan, Shiqi Wang, and Alex C Kot. Domain generalization with adversarial feature learning. In *CVPR*, 2018a.
- Ya Li, Xinmei Tian, Mingming Gong, Yajing Liu, Tongliang Liu, Kun Zhang, and Dacheng Tao. Deep domain generalization via conditional invariant adversarial networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 624–639, 2018b.
- Weixin Liang and James Zou. Metashift: A dataset of datasets for evaluating contextual distribution shifts and training conflicts. In *International Conference on Learning Representations*, 2022.
- Aengus Lynch, Gbètondji JS Dovonon, Jean Kaddour, and Ricardo Silva. Spawrious: A benchmark for fine control of spurious correlation biases. *arXiv preprint arXiv:2303.05470*, 2023.
- John P Miller, Rohan Taori, Aditi Raghunathan, Shiori Sagawa, Pang Wei Koh, Vaishaal Shankar, Percy Liang, Yair Carmon, and Ludwig Schmidt. Accuracy on the line: on the strong correlation between out-of-distribution and in-distribution generalization. In *Proceedings of the 38th International Conference on Machine Learning*, 2021.
- Hyeonseob Nam, HyunJae Lee, Jongchan Park, Wonjun Yoon, and Donggeun Yoo. Reducing domain gap by reducing style bias. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8690–8699, 2021.
- OpenAI. Gpt-4 technical report. Technical report, OpenAI, 2023.
- Giambattista Parascandolo, Alexander Neitz, ANTONIO ORVIETO, Luigi Gresele, and Bernhard Schölkopf. Learning explanations that are hard to vary. In *ICLR*, 2021.
- Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment matching for multi-source domain adaptation. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1406–1415, 2019.
- Mohammad Pezeshki, Oumar Kaba, Yoshua Bengio, Aaron C Courville, Doina Precup, and Guillaume Lajoie. Gradient starvation: A learning proclivity in neural networks. *Advances in Neural Information Processing Systems*, 34:1256–1272, 2021.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *Proceedings of the 38th International Conference on Machine Learning*, 2021.
- Yangjun Ruan, Yann Dubois, and Chris J Maddison. Optimal representations for covariate shift. *arXiv preprint arXiv:2201.00057*, 2021.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115:211–252, 2015.
- Shiori Sagawa, Pang Wei Koh, Tatsunori B. Hashimoto, and Percy Liang. Distributionally robust neural networks. In *ICLR*, 2020.

- Lukas Schott, Julius Von Kügelgen, Frederik Träuble, Peter Vincent Gehler, Chris Russell, Matthias Bethge, Bernhard Schölkopf, Francesco Locatello, and Wieland Brendel. Visual representation learning does not generalize strongly within the same domain. In *International Conference on Learning Representations*, 2022.
- Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pp. 618–626, 2017.
- Soroosh Shahtalebi, Jean-Christophe Gagnon-Audet, Touraj Laleh, Mojtaba Faramarzi, Kartik Ahuja, and Irina Rish. Sand-mask: An enhanced gradient masking strategy for the discovery of invariances in domain generalization. *arXiv preprint arXiv:2106.02266*, 2021.
- Zheyang Shen, Jiashuo Liu, Yue He, Xingxuan Zhang, Renzhe Xu, Han Yu, and Peng Cui. Towards out-of-distribution generalization: A survey. *arXiv:2108.13624*, 2021.
- Yuge Shi, Jeffrey Seely, Philip HS Torr, N Siddharth, Awni Hannun, Nicolas Usunier, and Gabriel Synnaeve. Gradient matching for domain generalization. *arXiv:2104.09937*, 2021.
- Yuge Shi, Imant Daunhawer, Julia E Vogt, Philip Torr, and Amartya Sanyal. How robust is unsupervised representation learning to distribution shift? In *The Eleventh International Conference on Learning Representations*, 2022.
- David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without human knowledge. *nature*, 550(7676):354–359, 2017.
- Baochen Sun, Jiashi Feng, and Kate Saenko. Correlation alignment for unsupervised domain adaptation. *Domain adaptation in computer vision applications*, pp. 153–171, 2017.
- Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In *ICLR*, 2014.
- Rohan Taori, Achal Dave, Vaishaal Shankar, Nicholas Carlini, Benjamin Recht, and Ludwig Schmidt. Measuring robustness to natural distribution shifts in image classification. *Advances in Neural Information Processing Systems*, 33:18583–18599, 2020.
- Vladimir Vapnik. *Statistical Learning Theory*. Wiley, 1998.
- Jindong Wang, Cuiling Lan, Chang Liu, Yidong Ouyang, Tao Qin, Wang Lu, Yiqiang Chen, Wenjun Zeng, and Philip Yu. Generalizing to unseen domains: A survey on domain generalization. *IEEE Transactions on Knowledge and Data Engineering*, 2022.
- Jindong Wang, Xixu Hu, Wenxin Hou, Hao Chen, Runkai Zheng, Yidong Wang, Linyi Yang, Haojun Huang, Wei Ye, Xiubo Geng, et al. On the robustness of chatgpt: An adversarial and out-of-distribution perspective. *arXiv preprint arXiv:2302.12095*, 2023.
- Olivia Wiles, Sven Gowal, Florian Stimberg, Sylvestre-Alvise Rebuffi, Ira Ktena, Krishnamurthy Dj Dvijotham, and Ali Taylan Cemgil. A fine-grained analysis on distribution shift. In *ICLR*, 2022.
- Mitchell Wortsman, Gabriel Ilharco, Jong Wook Kim, Mike Li, Simon Kornblith, Rebecca Roelofs, Raphael Gontijo Lopes, Hannaneh Hajishirzi, Ali Farhadi, Hongseok Namkoong, et al. Robust fine-tuning of zero-shot models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7959–7971, 2022.
- Yilun Xu and Tommi Jaakkola. Learning representations that support robust transfer of predictors. *arXiv preprint arXiv:2110.09940*, 2021.
- Shen Yan, Huan Song, Nanxiang Li, Lincan Zou, and Liu Ren. Improve unsupervised domain adaptation with mixup training. *arXiv:2001.00677*, 2020.

- Nanyang Ye, Kaican Li, Haoyue Bai, Runpeng Yu, Lanqing Hong, Fengwei Zhou, Zhenguo Li, and Jun Zhu. Ood-bench: Quantifying and understanding two dimensions of out-of-distribution generalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7947–7958, 2022.
- Guojun Zhang, Han Zhao, Yaoliang Yu, and Pascal Poupart. Quantifying and improving transferability in domain generalization. *Advances in Neural Information Processing Systems*, 34: 10957–10970, 2021a.
- Marvin Zhang, Henrik Marklund, Nikita Dhawan, Abhishek Gupta, Sergey Levine, and Chelsea Finn. Adaptive risk minimization: A meta-learning approach for tackling group distribution shift. *arXiv:2007.02931*, 2020.
- Marvin Zhang, Henrik Marklund, Nikita Dhawan, Abhishek Gupta, Sergey Levine, and Chelsea Finn. Adaptive risk minimization: Learning to adapt to domain shift. *Advances in Neural Information Processing Systems*, 34:23664–23678, 2021b.
- Bingchen Zhao, Shaozuo Yu, Wufei Ma, Mingxin Yu, Shenxiao Mei, Angtian Wang, Ju He, Alan Yuille, and Adam Kortylewski. Ood-cv: a benchmark for robustness to out-of-distribution shifts of individual nuisances in natural images. In *European Conference on Computer Vision*, pp. 163–180. Springer, 2022.
- Kaiyang Zhou, Ziwei Liu, Yu Qiao, Tao Xiang, and Chen Change Loy. Domain generalization: A survey. *arXiv:2103.02503*, 2021.

A ADDITIONAL INFORMATION ABOUT EXPERIMENT SETUP

A.1 DATASETS

Random examples drawn from each domain of the datasets we used (except NOISYMNIST which is shown in Figure 2) are shown in Figure 7-10. The order of the examples is arranged according to the degree of the distribution shift from low to high.



Figure 7: Examples of ROTATEDMNIST



Figure 8: Examples of LOWLIGHTCIFAR10



Figure 9: Examples of NOISYIMAGENET15



Figure 10: Examples of LR-IMAGENET15

For NOISYMNIST and ROTATEDMNIST, 60,000 images were divided into distinct training domains. For instance, in scenarios involving two training domains, each domain would encompass 30,000 images. Within the training domains, 20% of the data is allocated for in-distribution validation, aiding model calibration and selection. Every test domain of each altered dataset consists of 10,000 images, constructed using the same set of original images.

For NOISYIMAGENET15 and LR-IMAGENET15, we use the images in the training split of ImageNet to construct the training domains and the images in the validation split of ImageNet to construct the test domains. Similarly, the training domains divide the total 15,000 images in the training split. The test domains are constructed using the same set of original images, which consist of 750 images in total.

The 15 categories of birds we used in NOISYIMAGENET15 and LR-IMAGENET15, which correspond to indices 10 to 24 of the 1,000 categories of ImageNet, are “brambling, *Fringilla montifringilla*”, “goldfinch, *Carduelis carduelis*”, “house finch, linnnet, *Carpodacus mexicanus*”, “junco, snowbird”, “indigo bunting, indigo finch, indigo bird, *Passerina cyanea*”, “robin, American robin, *Turdus migratorius*”, “bulbul”, “jay”, “magpie”, “chickadee”, “water ouzel, dipper”, “kite”, “bald eagle, American eagle, *Haliaeetus leucocephalus*”, “vulture”, and “great grey owl, great gray owl, *Strix nebulosa*”.

In addition to the above datasets studied in the paper, we have also conducted preliminary experiments on another dataset called IMPULSENISEMNIST. The dataset is constructed by gradually adding impulse noise to MNIST. Below are some of the examples of this dataset. The experiment results on this dataset are given in Appendix B.2.



Figure 11: Examples of IMPULSENOISEMNIST

A.2 ALGORITHMS

Here is the full list of domain generalization algorithms we used in this study:

- Invariant Risk Minimization (**IRM**, [Arjovsky et al. \(2019\)](#))
- Group Distributionally Robust Optimization (**GroupDRO**, [Sagawa et al. \(2020\)](#))
- Interdomain Mixup (**Mixup**, [Yan et al. \(2020\)](#))
- Marginal Transfer Learning (**MTL**, [Blanchard et al. \(2011\)](#))
- Maximum Mean Discrepancy (**MMD**, [Li et al. \(2018a\)](#))
- Deep CORAL (**CORAL**, [Sun et al. \(2017\)](#))
- Domain Adversarial Neural Network (**DANN**, [Ganin et al. \(2016\)](#))
- Conditional Domain Adversarial Neural Network (**CDANN**, [Li et al. \(2018b\)](#))
- Style Agnostic Networks (**SagNet**, [Nam et al. \(2021\)](#))
- Adaptive Risk Minimization (**ARM**, [Zhang et al. \(2021b\)](#))
- Variance Risk Extrapolation (**VREx**, [Krueger et al. \(2020\)](#))
- Representation Self-Challenging (**RSC**, [Huang et al. \(2020\)](#))
- Spectral Decoupling (**SD**, [Pezeshki et al. \(2021\)](#))
- Learning Explanations that are Hard to Vary (**AND-Mask**, [Parascandolo et al. \(2021\)](#))
- Smoothed-AND mask (**SAND-mask**, [Shahtalebi et al. \(2021\)](#))
- Out-of-Distribution Generalization with Maximal Invariant Predictor (**IGA**, [Koyama & Yamaguchi \(2020\)](#))
- Gradient Matching for Domain Generalization (**Fish**, [Shi et al. \(2021\)](#))
- Self-supervised Contrastive Regularization (**SelfReg**, [Kim et al. \(2021\)](#))
- Learning Representations that Support Robust Transfer of Predictors (**TRM**, [Xu & Jaakkola \(2021\)](#))
- Invariance Principle Meets Information Bottleneck for Out-of-Distribution Generalization (**IB-ERM** & **IB-IRM**, [Ahuja et al. \(2021\)](#))
- Optimal Representations for Covariate Shift (**CAD** & **CondCAD**, [Ruan et al. \(2021\)](#))
- Quantifying and Improving Transferability in Domain Generalization (**Transfer**, [Zhang et al. \(2021a\)](#))
- Invariant Causal Mechanisms through Distribution Matching (**CausIRL** with CORAL or MMD, [Chevalley et al. \(2022\)](#))
- Empirical Quantile Risk Minimization (**EQRM**, [Eastwood et al. \(2022\)](#))

We use the DomainBed ([Gulrajani & Lopez-Paz, 2021](#)) implementation for all the above algorithms.

A.3 ADDITIONAL IMPLEMENTATION DETAILS

Except for learning rate, batch size, weight decay, and dropout, the search of other hyperparameters follows that of DomainBed ([Gulrajani & Lopez-Paz, 2021](#)). For experiments utilizing the ResNet-50 architecture, the original MNIST digits were resized to a resolution of 224x224 pixels. Subsequent normalization was performed using the mean and standard deviation inherent to the MNIST dataset to preserve data integrity and distribution.

B ADDITIONAL EXPERIMENT RESULTS

B.1 EFFECTS OF INCREASING THE NUMBER OF TRAINING DOMAINS

In Figure 3, we have shown the best-performing models (trained on \mathcal{D}_0 and \mathcal{D}_1) at each degree of NOISYMNIST. In Figure 12 and Figure 13 below, we show more results on this dataset, with models trained on wider ranges of domains. The results show that more training domains help. As more training domains are added, the gaps between the best-performing models gradually decrease. Nevertheless, the gaps seem to be closing at a slow rate. The discrepancy between the best-performing models is still largely present. In addition, the advantage of DG methods over ERM becomes less pronounced as the number of training domains increases.

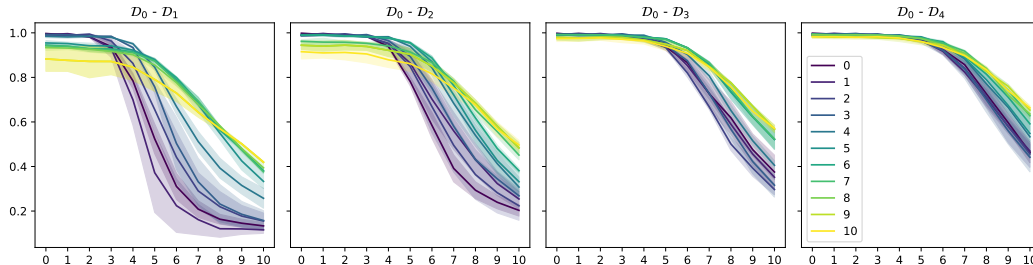


Figure 12: Performance of the best-performing models at each degree of NOISYMNIST. The label of the curves denotes the domain on which the models perform best. The title of each subplot denotes the training domains (e.g., “ $\mathcal{D}_0 - \mathcal{D}_2$ ” means that the models are trained on $\mathcal{D}_0, \mathcal{D}_1, \mathcal{D}_2$). The results are averaged over the top 3 models of all algorithms at each degree.

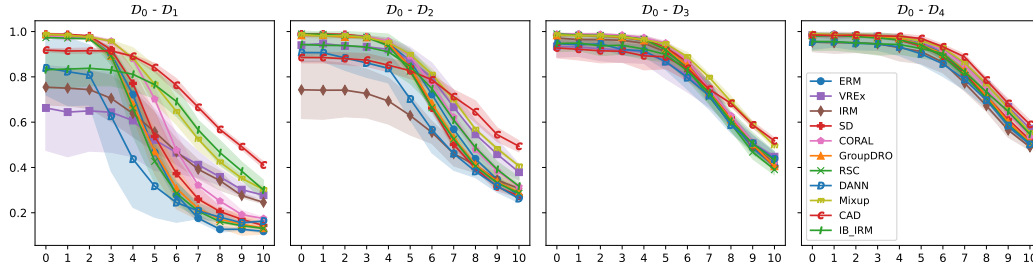


Figure 13: Performance of ERM and representative domain generalization algorithms on NOISYMNIST. The title of each subplot denotes the training domains (e.g., “ $\mathcal{D}_0 - \mathcal{D}_2$ ” means that the models are trained on $\mathcal{D}_0, \mathcal{D}_1, \mathcal{D}_2$). The results are averaged over the top 3 models of each algorithm.

B.2 EXPERIMENT RESULTS ON IMPULSENOISEMNIST

In Figure 14, we show the experiment results on the IMPULSENOISEMNIST dataset. The results on the generalization from milder to stronger shifts are largely consistent with the results in Figure 3. As for the results on the generalization from stronger to milder shifts, the pattern looks like that of LOWLIGHTCIFAR10. In particular, when the shift in the training data is not very strong (e.g., \mathcal{D}_0 and \mathcal{D}_5), the generalization pattern looks just like that of NoisyMNIST—robustness against milder shifts is not affected. On the other hand, when the shift in the training data is very strong (e.g., \mathcal{D}_0 and \mathcal{D}_{10}), the pattern looks like that of RotatedMNIST—robustness against milder shifts is significantly weakened.

B.3 STRONG-TO-MILD GENERALIZATION PERFORMANCE OF DG ALGORITHMS

In this section, we show the performance of DG algorithms for the hardest generalization case in Figure. 5. From Figure 15 to Figure 17, we can see that most DG algorithms are helpful to the

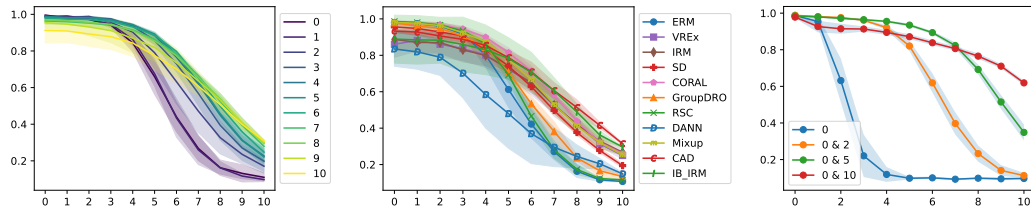


Figure 14: Results on IMPULSENNOISEMNIST (Simple CNN). **(Left)** Performance of the best-performing models at each shift degree. The label of the curves denotes the domain on which the models perform best; **(Middle)** Performance of ERM and representative domain generalization algorithms on the two datasets; **(Right)** Performance of ERM models trained on domains under different shift degrees. The label of the curves denotes the indices of the training domains, e.g., “0 & 2” means that the models are trained on \mathcal{D}_0 and \mathcal{D}_2 of the corresponding dataset. The implementation details of these experiments follow those of NOISYMNIIST.

generalization from stronger shifts to milder shifts in the case of ROTATEDMNIST, although only to a limited extent. On the other two datasets, only some of the DG algorithms are able to improve the generalization performance by a very small margin.

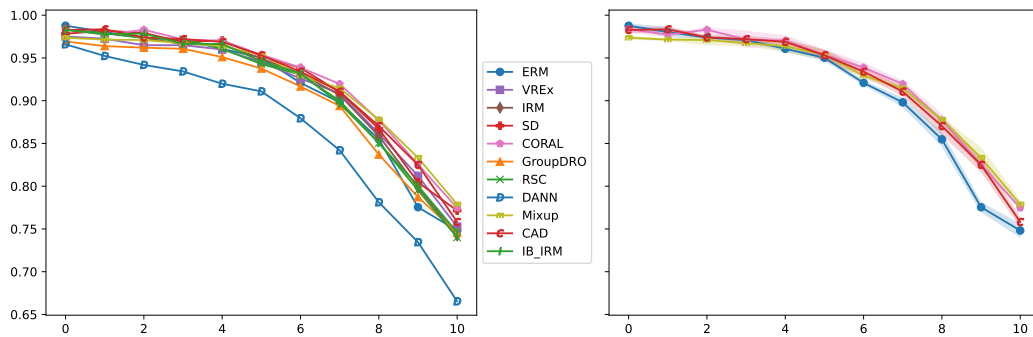


Figure 15: Average accuracy of the top-3 models of ERM and various DG algorithms trained on \mathcal{D}_0 and \mathcal{D}_{10} of NOISYMNIIST. The models are selected via training-domain validation. Error bars are omitted in the left sub-figure for clarity. The best-performing DG algorithms are compared with ERM in the right sub-figure. The results are averaged over 3 runs.

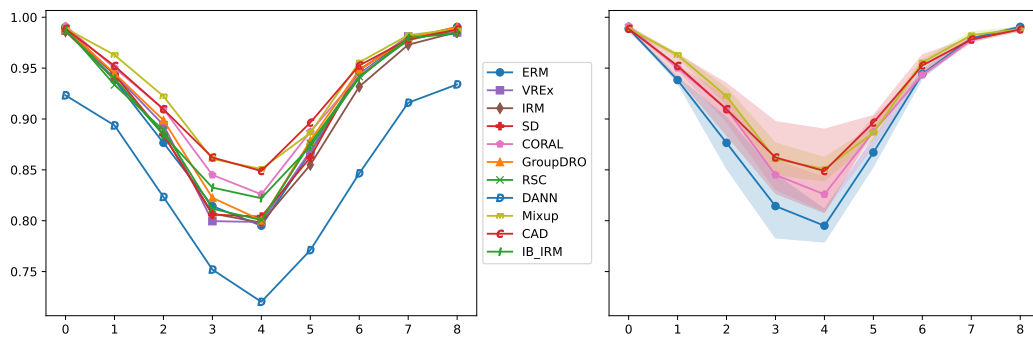


Figure 16: Average accuracy of the top-3 models of ERM and various DG algorithms trained on \mathcal{D}_0 and \mathcal{D}_8 of ROTATEDMNIST. The models are selected via training-domain validation. Error bars are omitted in the left sub-figure for clarity. The best-performing DG algorithms are compared with ERM in the right sub-figure. The results are averaged over 3 runs.

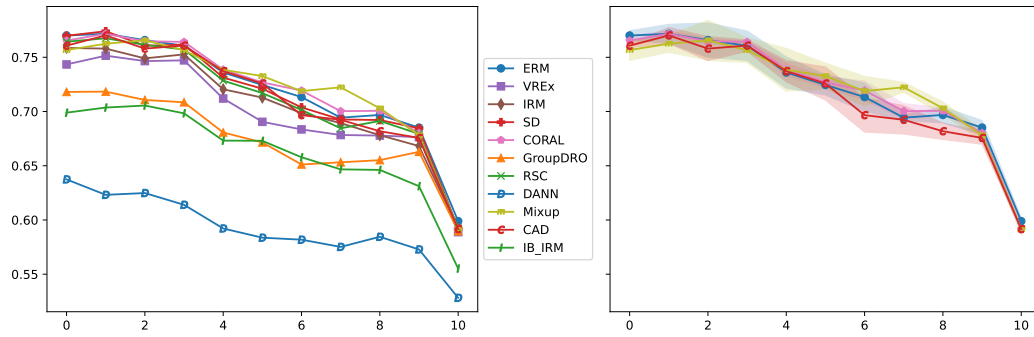


Figure 17: Average accuracy of the top-3 models of ERM and various DG algorithms trained on \mathcal{D}_0 and \mathcal{D}_{10} of LOWLIGHTCIFAR10. The models are selected via training-domain validation. Error bars are omitted in the left sub-figure for clarity. The best-performing DG algorithms are compared with ERM in the right sub-figure. The results are averaged over 3 runs.

B.4 FULL NUMERICAL RESULTS OF THE EXPERIMENTS CONDUCTED IN SECTION 4

From Table 2 to Table 19, we show the full numerical results of the experiments we conducted in Section 4. The first column of each table shows the domain on which the accuracy is used to select the top models.

Algorithm	\mathcal{D}_0	\mathcal{D}_1	\mathcal{D}_2	\mathcal{D}_3	\mathcal{D}_4	\mathcal{D}_5	\mathcal{D}_6	\mathcal{D}_7	\mathcal{D}_8	\mathcal{D}_9	\mathcal{D}_{10}	
\mathcal{D}_8	ERM	93.96±0.57	93.89±0.69	94.55±1.80	76.12±9.31	51.92±10.73	34.90±5.89	25.72±3.10	19.13±1.45	16.65±0.55	14.53±1.22	14.08±1.19
	VREx	68.05±39.25	68.24±39.24	66.77±37.74	49.25±26.67	27.30±12.21	20.68±7.95	18.20±4.25	15.98±2.05	15.23±0.16	12.92±0.28	11.45±0.70
	IRM	41.32±8.56	42.73±9.53	54.90±9.10	49.45±5.47	43.55±5.23	37.28±4.48	33.27±4.14	28.63±1.86	25.07±1.83	22.15±1.17	20.40±1.21
	SD	96.73±1.97	96.94±1.79	96.97±1.64	85.37±3.78	55.72±9.91	36.37±6.48	27.62±3.43	20.22±0.98	18.20±0.14	15.47±0.15	14.15±0.25
	CORAL	97.69±0.39	97.43±0.34	96.57±0.64	75.20±4.90	45.62±3.97	29.60±4.40	23.87±2.71	18.08±1.69	16.68±0.83	14.12±1.18	14.17±0.45
	GroupDRO	96.16±2.12	96.34±2.15	94.83±2.94	70.40±10.43	45.82±11.71	31.45±6.90	24.10±2.27	18.45±1.46	16.55±0.37	14.87±1.09	13.80±1.33
	RSC	96.63±2.31	96.41±2.13	95.53±3.64	83.10±5.90	57.03±6.74	38.67±4.21	29.10±2.41	20.52±1.19	17.13±0.27	13.88±0.90	12.88±0.39
	DANN	90.71±3.94	90.24±4.45	91.12±3.27	74.17±0.95	50.58±4.66	34.43±3.59	24.62±1.83	19.10±0.60	17.08±0.47	13.95±0.58	14.23±1.47
	Mixup	92.65±3.11	93.67±2.93	73.47±20.07	54.52±21.52	40.12±14.61	27.37±7.39	21.68±2.32	18.58±0.93	17.43±0.29	15.97±1.29	15.70±1.64
	CAD	93.47±6.63	93.63±6.24	94.08±6.29	83.62±7.76	57.38±3.27	34.82±3.47	25.82±3.85	20.18±3.82	18.48±2.02	15.82±1.92	14.70±1.32
IB-IRM	63.79±17.64	64.18±16.20	70.88±15.05	57.33±15.06	38.02±6.78	26.87±3.46	22.50±2.58	19.48±2.04	18.02±1.23	16.20±0.42	15.57±1.18	
\mathcal{D}_9	ERM	93.77±0.69	93.52±0.66	92.75±1.20	67.17±6.77	43.17±8.51	29.28±5.51	22.17±3.40	17.82±2.30	16.07±1.27	15.00±0.74	14.23±1.00
	VREx	51.79±32.78	52.77±31.84	53.92±31.61	47.82±25.99	31.95±11.56	22.70±2.55	18.27±0.93	16.28±0.63	14.02±0.25	14.47±0.91	13.33±0.66
	IRM	41.32±8.56	42.73±9.53	54.90±9.10	49.45±5.47	43.55±5.23	37.28±4.48	33.27±4.14	28.63±1.86	25.07±1.83	22.15±1.17	20.40±1.21
	SD	96.66±1.92	96.94±1.78	96.95±1.63	84.97±3.22	56.30±10.54	36.18±6.26	27.53±3.33	20.57±1.47	18.07±0.17	15.88±0.44	14.15±0.25
	CORAL	96.98±1.02	96.66±1.34	92.43±4.34	63.47±8.20	36.53±9.47	25.22±7.67	20.97±4.78	17.10±2.46	15.75±1.46	15.20±0.52	14.52±0.51
	GroupDRO	96.43±2.33	96.49±2.26	94.03±2.83	72.13±9.73	49.53±8.92	32.82±6.10	23.88±2.42	18.45±1.46	16.28±0.53	15.52±0.35	14.07±1.12
	RSC	98.04±0.23	97.53±0.39	97.93±0.45	86.55±2.30	53.83±3.57	32.85±0.82	23.63±3.04	18.20±1.94	16.02±1.18	14.42±0.53	13.35±0.98
	DANN	91.91±4.32	91.54±4.72	92.98±3.49	76.97±3.03	53.08±5.22	35.58±3.91	25.57±1.57	18.80±0.52	17.08±0.47	14.47±0.55	15.27±0.02
	Mixup	93.12±3.45	93.12±2.45	60.37±21.38	42.02±26.76	32.47±17.76	23.30±8.86	19.67±2.65	18.10±1.06	17.18±0.47	16.70±0.60	15.98±1.53
	CAD	91.11±5.83	91.44±5.54	89.23±6.90	72.55±14.35	48.60±9.63	31.42±6.51	25.45±4.27	20.60±3.37	18.25±2.33	16.08±1.71	14.22±1.88
IB-IRM	50.49±18.28	50.21±17.76	57.20±19.54	46.30±17.34	35.55±11.15	27.73±6.54	23.12±3.63	19.72±1.60	17.57±1.52	16.57±0.77	16.03±0.34	
\mathcal{D}_{10}	ERM	95.04±1.51	94.83±1.74	94.77±2.08	73.92±6.85	45.35±7.38	29.13±5.65	22.23±3.32	18.08±2.02	16.23±1.05	14.92±0.81	14.70±0.53
	VREx	71.82±33.71	72.01±32.79	72.85±32.82	61.10±25.19	36.80±10.57	23.50±1.53	17.93±1.40	15.70±1.41	14.07±0.18	14.03±1.29	13.53±0.47
	IRM	41.32±8.56	42.73±9.53	54.90±9.10	49.45±5.47	43.55±5.23	37.28±4.48	33.27±4.14	28.63±1.86	25.07±1.83	22.15±1.17	20.40±1.21
	SD	96.84±2.06	96.91±1.76	96.73±1.58	84.48±4.47	54.75±11.21	35.23±7.93	26.00±5.51	19.43±1.84	17.13±1.65	15.33±0.19	14.38±0.10
	CORAL	97.08±1.05	96.87±1.41	92.73±4.58	61.57±5.54	33.35±5.02	21.73±3.28	18.80±1.98	16.00±1.06	15.30±0.83	14.48±0.71	14.82±0.27
	GroupDRO	97.69±0.14	97.86±0.14	97.20±0.36	82.75±2.71	58.05±3.52	38.75±2.79	26.97±2.38	20.18±1.25	16.48±0.43	15.22±0.60	14.82±0.64
	RSC	98.17±0.07	97.92±0.29	97.85±0.35	84.32±3.48	50.62±3.83	29.58±3.89	21.40±3.92	16.73±2.03	15.20±1.50	14.32±0.64	13.78±0.59
	DANN	91.91±4.32	91.54±4.72	92.98±3.49	76.97±3.03	53.08±5.22	35.58±3.91	25.57±1.57	18.80±0.52	17.08±0.47	14.47±0.55	15.27±0.02
	Mixup	96.03±0.77	95.80±1.34	54.08±12.49	29.55±9.47	22.45±3.83	18.52±2.16	18.45±1.00	17.50±0.21	17.00±0.57	16.42±0.94	16.50±1.06
	CAD	93.47±6.63	93.63±6.24	94.08±6.29	83.62±7.76	57.38±3.27	34.82±3.47	25.82±3.85	20.18±3.82	18.48±2.02	15.82±1.92	14.70±1.32
IB-IRM	55.73±25.00	55.51±24.36	61.58±24.72	49.53±21.59	34.65±9.97	26.07±4.41	22.42±2.64	19.78±1.69	17.63±1.61	16.28±0.38	16.12±0.42	

Table 7: Average accuracy of top-3 models of each algorithm at each shift degree of NOISYMNIST (EfficientNet-b0). This table shows the results on shift degrees from \mathcal{D}_8 to \mathcal{D}_{10} .

Table with columns for Algorithm, D0, D1, D2, D3, D4, D5, D6, D7, D8, D9, D10. Rows are grouped by algorithm (ERM, VREx, IRM, SD, GroupDRO, RSC, Mixup, CORAL, MMD, DANN, MTL, SagNet, ARM, ANDMask, SelfReg, Fish, IB-ERM, IB-IRM, CAD, CondCAD, EQRM) and shift degree (D4, D5, D6, D7).

Table 9: Average accuracy of top-3 models of each algorithm at each shift degree of RESISYMNIST (ResNet-50). This table shows the results on shift degrees from D4 to D7.

Algorithm	\mathcal{D}_0	\mathcal{D}_1	\mathcal{D}_2	\mathcal{D}_3	\mathcal{D}_4	\mathcal{D}_5	\mathcal{D}_6	\mathcal{D}_7	\mathcal{D}_8	\mathcal{D}_9	\mathcal{D}_{10}	
\mathcal{D}_8	ERM	75.98±2.23	75.15±3.41	18.30±0.76	19.60±0.86	19.13±1.11	19.80±0.39	19.43±0.37	19.48±0.87	19.12±0.41	18.57±0.68	17.42±0.52
	VREx	60.56±1.77	59.70±2.61	15.73±2.86	16.55±1.76	16.73±1.64	18.22±1.74	18.93±0.49	19.40±0.51	18.03±0.56	16.18±0.96	15.65±1.75
	IRM	44.99±26.43	45.65±26.30	16.00±3.32	16.40±2.18	16.45±3.28	17.27±2.65	17.53±1.10	17.18±0.92	17.20±0.73	15.07±0.74	14.92±1.10
	SD	73.58±1.90	74.13±2.25	16.60±0.63	18.38±0.84	18.20±0.78	19.45±1.06	19.82±2.09	20.72±1.25	19.17±0.10	17.60±0.14	17.53±1.17
	CORAL	81.36±3.34	81.50±3.33	19.02±0.55	18.68±0.84	17.95±0.44	19.32±0.24	19.22±0.93	18.78±0.95	18.48±0.58	16.57±1.78	16.37±2.22
	GroupDRO	74.78±9.15	75.56±8.89	17.23±0.80	18.60±0.59	18.60±1.20	19.03±1.05	18.60±1.45	18.23±1.09	18.97±0.37	16.78±0.93	16.50±0.71
	RSC	72.15±10.29	72.61±10.78	15.82±0.22	15.83±0.61	15.10±0.63	16.28±0.31	16.73±0.66	16.18±0.41	16.77±0.55	15.70±0.63	15.68±0.48
	DANN	50.20±15.53	51.31±14.61	16.98±3.08	17.58±2.84	17.08±2.29	18.52±2.72	17.80±1.74	16.82±1.74	16.68±0.18	15.28±0.73	13.85±0.56
	Mixup	67.31±8.62	67.48±8.63	17.67±1.24	18.18±1.48	17.33±1.36	18.67±0.97	18.77±0.72	18.67±1.18	18.37±0.53	16.27±0.93	15.77±1.53
	CAD	66.80±16.15	66.62±16.00	18.40±1.82	18.78±0.83	17.82±0.59	19.47±0.83	19.45±0.90	19.23±0.57	19.50±0.39	18.32±0.73	18.22±0.26
	IB-IRM	50.17±25.97	50.22±25.79	17.05±1.70	17.30±1.36	16.87±1.05	17.85±1.04	17.20±0.29	16.33±0.17	16.57±0.57	15.13±0.78	14.72±1.05
	\mathcal{D}_9	ERM	79.09±2.83	78.59±3.20	17.77±1.00	19.02±1.24	18.50±1.35	19.00±1.02	19.52±0.35	19.18±1.10	18.82±0.74	18.80±0.46
VREx		66.77±10.37	66.40±10.73	16.65±1.62	16.75±1.53	16.92±1.38	18.08±1.90	18.30±1.38	18.97±1.08	17.58±1.09	16.50±0.59	15.80±1.63
IRM		46.58±23.15	46.09±22.42	13.65±1.90	15.18±1.62	15.18±2.30	15.62±1.39	17.08±0.99	16.25±0.46	16.33±0.12	15.70±1.27	15.50±0.62
SD		61.63±0.75	63.10±0.54	14.37±0.52	15.83±0.71	15.83±0.85	16.57±1.47	18.60±1.04	19.48±1.23	18.18±0.59	18.77±0.33	18.32±0.57
CORAL		70.51±5.86	70.55±5.71	18.17±1.86	19.58±1.37	18.30±0.98	19.22±0.29	18.90±1.11	17.88±1.57	17.75±1.17	17.12±1.37	16.50±2.05
GroupDRO		69.58±9.37	69.59±9.38	16.25±0.92	17.05±1.32	17.43±1.38	17.58±1.02	17.23±0.88	17.20±0.37	17.77±0.52	17.20±0.57	16.78±0.64
RSC		71.41±9.63	71.45±9.78	16.22±0.60	16.33±0.92	15.97±0.88	17.25±1.11	17.80±1.07	16.57±0.41	16.55±0.43	15.78±0.74	15.35±0.78
DANN		64.61±17.95	65.02±17.14	17.70±2.23	18.52±2.04	17.33±2.09	18.20±2.52	17.93±1.43	17.10±1.25	16.43±0.49	16.58±1.22	15.48±1.29
Mixup		66.76±10.16	66.93±10.04	17.42±1.02	18.25±1.24	17.27±1.38	18.67±0.72	18.57±0.05	19.33±0.35	17.93±0.81	17.35±0.32	16.47±0.06
CAD		66.36±15.72	66.43±15.80	17.65±0.99	18.37±0.27	18.12±0.94	19.32±0.62	19.68±0.73	19.08±0.65	19.32±0.63	18.55±0.40	17.85±0.78
IB-IRM		50.43±26.27	50.35±25.94	15.97±0.50	16.27±0.42	15.78±0.94	16.47±1.02	16.70±0.94	15.37±1.33	15.87±1.45	15.17±0.73	14.27±1.68
\mathcal{D}_{10}		ERM	72.55±7.32	72.05±7.16	16.42±3.13	17.93±2.32	17.57±1.88	18.50±1.66	19.05±0.14	18.48±0.84	18.67±0.66	17.60±0.28
	VREx	59.83±0.99	59.85±2.70	17.33±0.83	17.85±1.00	17.85±0.27	18.95±0.99	19.10±0.27	18.38±1.88	17.77±0.86	16.08±1.09	15.87±1.58
	IRM	43.40±20.38	43.37±20.15	13.82±2.11	15.87±1.79	15.05±2.24	16.03±1.47	17.32±0.99	16.37±0.63	16.67±0.39	15.53±1.08	15.83±0.18
	SD	68.13±9.91	69.50±9.48	14.95±1.23	16.08±0.89	16.03±1.05	16.12±1.35	18.23±0.97	18.57±1.15	18.35±0.62	18.35±0.55	18.78±0.38
	CORAL	81.49±3.43	81.30±3.18	17.47±1.37	18.18±1.10	17.13±1.02	18.25±0.66	18.35±1.50	17.98±1.49	17.80±1.03	16.55±1.83	17.17±1.36
	GroupDRO	68.73±9.89	69.32±9.55	16.37±1.08	17.40±1.66	17.30±1.20	17.93±1.47	17.88±1.79	17.70±1.02	18.28±0.69	17.17±0.60	17.17±0.25
	RSC	73.23±10.96	73.07±10.93	14.92±0.86	14.98±1.15	14.55±0.99	15.25±1.09	15.93±1.28	16.02±0.47	16.08±0.55	15.38±0.34	16.28±0.77
	DANN	66.21±18.60	66.33±17.66	17.97±2.60	18.52±2.04	17.20±1.91	17.80±1.95	17.55±0.90	16.78±0.80	16.28±0.45	16.12±1.48	15.58±1.21
	Mixup	60.08±13.67	61.42±12.98	16.22±2.02	17.02±1.98	15.92±2.08	17.27±2.48	17.75±1.95	17.97±1.93	17.43±0.71	16.43±0.24	17.18±0.06
	CAD	63.55±15.39	63.79±15.30	16.62±1.17	17.63±0.57	17.25±0.54	18.30±0.49	18.33±0.53	18.33±0.23	19.03±0.65	18.40±0.49	18.65±0.98
	IB-IRM	50.85±26.77	50.97±26.68	16.03±0.52	16.62±0.52	15.87±0.86	16.57±0.89	16.88±0.69	15.45±1.21	15.97±1.31	14.73±1.34	14.77±0.98

Table 19: Average accuracy of top-3 models of each algorithm at each shift degree of LOWLIGHT-CIFAR10 (ResNet-50). This table shows the results on shift degrees from \mathcal{D}_8 to \mathcal{D}_{10} .

B.5 COMPLETE RESULTS OF LINEAR PROBING

	\mathcal{D}_0	\mathcal{D}_1	\mathcal{D}_2	\mathcal{D}_3	\mathcal{D}_4	\mathcal{D}_5	\mathcal{D}_6	\mathcal{D}_7	\mathcal{D}_8	\mathcal{D}_9	\mathcal{D}_{10}
IN ₀	98.4±0.1	39.8±4.2	24.8±4.7	16.7±3.2	12.4±2.1	10.9±2.0	10.3±1.5	10.1±1.4	9.8±1.1	9.6±0.8	9.5±0.6
IN ₁	97.9±0.1	97.3±0.2	90.2±0.7	62.9±2.0	33.1±1.7	20.1±1.9	15.4±1.7	13.3±1.1	11.7±1.1	10.8±0.9	10.5±1.0
IN ₂	97.6±0.2	97.6±0.4	95.8±0.4	88.2±0.5	65.4±0.8	39.9±1.8	25.6±2.6	19.1±2.8	15.7±2.1	13.5±1.3	12.2±1.0
IN ₃	97.6±0.2	97.2±0.1	95.9±0.1	92.6±0.8	81.8±0.4	57.9±1.3	36.2±1.6	23.5±1.7	17.4±1.9	14.3±1.6	12.7±1.1
IN ₄	96.6±0.3	96.6±0.5	96.1±0.6	93.5±0.6	87.4±1.0	74.4±0.6	54.4±1.4	37.3±2.0	26.5±1.7	20.4±1.1	16.9±0.9
CLIP ₀	98.1±0.1	55.4±4.1	35.9±4.4	26.4±3.3	21.1±2.7	18.1±1.4	15.8±0.4	14.1±0.4	13.0±0.4	12.3±0.5	12.1±0.2
CLIP ₁	97.9±0.2	95.9±0.2	86.3±0.7	63.0±2.1	41.2±2.4	26.2±1.2	18.3±1.0	14.5±1.0	12.8±0.8	11.8±0.8	11.4±0.7
CLIP ₂	97.4±0.2	96.2±0.2	92.7±0.5	81.0±0.6	58.3±1.7	35.9±1.7	22.9±0.6	16.9±0.4	14.0±0.3	12.8±0.6	11.8±0.4
CLIP ₃	97.0±0.3	95.7±0.4	93.0±0.7	85.9±0.6	70.6±0.2	45.9±1.4	24.0±2.1	15.0±1.6	11.7±0.8	10.8±0.5	10.3±0.2
CLIP ₄	96.1±0.4	95.7±0.5	92.9±0.1	87.1±0.7	76.9±1.0	59.3±0.5	38.9±0.6	24.6±1.0	17.0±1.0	13.6±0.6	12.2±0.4

Table 20: ResNet-50 results on NOISYMNIST.

	\mathcal{D}_0	\mathcal{D}_1	\mathcal{D}_2	\mathcal{D}_3	\mathcal{D}_4	\mathcal{D}_5	\mathcal{D}_6	\mathcal{D}_7	\mathcal{D}_8	\mathcal{D}_9	\mathcal{D}_{10}
IN ₀	98.4±0.3	69.7±1.4	51.7±1.7	36.1±1.6	25.6±1.5	17.6±0.8	13.5±0.4	11.8±0.4	11.6±0.7	11.1±0.5	10.8±0.6
IN ₁	98.1±0.2	97.6±0.2	94.2±0.2	82.2±1.6	58.4±4.7	35.2±4.7	23.0±2.3	18.0±1.0	15.1±0.5	13.9±0.4	13.1±0.7
IN ₂	98.1±0.1	97.7±0.4	96.3±0.3	92.2±0.2	78.5±0.6	52.0±1.1	30.0±1.5	19.3±1.7	14.5±1.0	12.5±0.4	12.1±0.5
IN ₃	98.0±0.5	97.4±0.1	96.5±0.3	94.7±0.5	85.7±0.4	64.0±1.2	39.4±2.2	24.3±1.8	17.6±1.5	14.6±1.5	13.4±1.3
IN ₄	97.4±0.6	96.6±0.5	96.7±0.2	94.3±0.4	89.1±0.6	75.2±0.6	53.0±2.1	32.4±3.0	21.1±2.8	15.8±2.3	13.3±1.6
CLIP ₀	98.8±0.1	89.7±1.9	71.2±3.7	48.2±2.3	31.2±1.3	21.7±1.0	16.3±1.0	13.7±1.0	12.3±1.1	11.6±1.1	11.2±1.2
CLIP ₁	98.5±0.2	97.8±0.2	91.5±1.0	69.3±4.2	45.0±5.6	29.9±5.5	22.1±4.1	17.8±3.1	15.4±2.0	14.2±1.9	12.9±1.3
CLIP ₂	98.7±0.1	98.0±0.4	95.6±0.4	86.3±0.6	65.7±1.3	43.1±0.6	29.7±0.4	22.6±0.6	18.8±0.5	16.8±0.1	15.8±0.3
CLIP ₃	98.3±0.2	97.5±0.1	96.0±0.1	90.4±0.6	77.2±0.4	56.7±0.6	38.8±0.2	27.6±0.5	21.8±0.4	18.1±0.5	16.3±0.5
CLIP ₄	98.0±0.6	97.4±0.3	95.8±0.6	91.0±0.4	81.3±0.8	67.1±0.7	50.1±0.9	36.6±0.3	27.6±0.5	22.3±0.4	19.1±0.3

Table 21: ViT-B/32 results on NOISYMNIST.

	\mathcal{D}_0	\mathcal{D}_1	\mathcal{D}_2	\mathcal{D}_3	\mathcal{D}_4	\mathcal{D}_5	\mathcal{D}_6	\mathcal{D}_7	\mathcal{D}_8	\mathcal{D}_9	\mathcal{D}_{10}
IN ₀	98.5±0.1	97.3±0.1	94.7±0.1	89.2±0.2	80.3±0.2	68.1±0.0	55.0±0.3	44.5±0.4	38.6±0.7	36.5±0.8	34.3±1.2
IN ₁	98.4±0.1	98.6±0.2	97.4±0.1	94.5±0.1	88.3±0.2	77.4±0.2	64.2±0.6	53.2±0.9	46.2±1.1	42.2±1.4	38.6±1.2
IN ₂	98.2±0.2	98.5±0.3	98.3±0.1	97.0±0.1	93.8±0.3	87.2±0.5	77.1±1.1	67.1±0.9	58.6±1.4	52.4±2.0	48.6±2.1
IN ₃	98.1±0.2	98.6±0.3	98.5±0.1	98.6±0.2	96.7±0.1	93.4±0.0	86.7±0.2	78.1±0.8	68.5±1.4	58.9±1.6	54.6±0.7
IN ₄	97.2±0.2	98.0±0.3	98.4±0.2	98.5±0.6	97.8±0.4	96.4±0.1	92.9±0.2	86.8±0.1	77.6±0.9	65.9±1.4	61.3±1.5
CLIP ₀	98.2±0.1	96.6±0.1	93.4±0.2	83.6±1.0	69.6±1.4	55.4±1.4	44.8±0.7	40.5±0.9	38.0±1.7	37.9±0.6	34.2±1.9
CLIP ₁	98.0±0.1	98.2±0.1	96.3±0.2	90.0±0.3	78.6±0.6	64.7±0.9	53.1±0.6	48.1±0.8	43.0±0.8	40.7±0.7	37.6±0.6
CLIP ₂	97.7±0.1	98.0±0.1	97.5±0.1	95.4±0.2	90.1±0.2	81.0±0.7	69.9±0.9	60.7±1.1	50.5±0.7	44.2±0.7	39.6±1.5
CLIP ₃	97.6±0.2	98.1±0.1	97.9±0.5	97.6±0.4	95.3±0.1	90.5±0.3	82.1±0.5	71.5±0.6	58.3±0.8	50.8±0.6	44.1±1.4
CLIP ₄	97.2±0.4	97.6±0.2	98.0±0.5	97.8±0.3	96.7±0.4	94.7±0.1	89.2±0.2	80.6±0.5	68.0±0.9	56.4±0.9	49.3±1.0

Table 22: ResNet-50 results on ROTATEDMNIST.

	\mathcal{D}_0	\mathcal{D}_1	\mathcal{D}_2	\mathcal{D}_3	\mathcal{D}_4	\mathcal{D}_5	\mathcal{D}_6	\mathcal{D}_7	\mathcal{D}_8	\mathcal{D}_9	\mathcal{D}_{10}
IN ₀	98.3±0.2	97.7±0.0	96.1±0.1	90.3±0.3	80.8±0.5	68.4±0.7	55.7±0.8	46.7±0.1	41.2±0.3	39.8±0.6	38.9±0.4
IN ₁	98.3±0.0	98.6±0.1	97.8±0.1	95.0±0.2	89.2±0.4	79.5±0.5	67.7±0.6	58.2±0.9	51.1±0.3	45.9±0.9	45.0±0.9
IN ₂	98.2±0.1	98.7±0.3	98.7±0.2	97.4±0.1	94.4±0.2	88.3±0.5	78.0±1.6	66.6±1.9	57.7±1.6	50.8±1.5	48.9±0.7
IN ₃	97.9±0.2	98.4±0.3	98.7±0.4	98.6±0.1	96.9±0.1	94.1±0.3	87.5±0.3	77.7±0.2	67.7±0.6	59.7±1.0	56.6±0.8
IN ₄	97.4±0.7	98.2±0.5	99.0±0.2	98.5±0.3	98.4±0.4	96.7±0.1	92.8±0.1	85.4±0.2	76.7±0.3	66.6±0.1	62.0±0.3
CLIP ₀	98.8±0.1	97.8±0.1	95.1±0.1	87.5±0.4	73.9±0.7	58.0±0.6	44.5±0.5	37.6±0.4	33.3±0.1	32.1±0.5	33.3±0.2
CLIP ₁	98.6±0.2	98.9±0.1	97.7±0.1	93.3±0.3	84.2±0.6	71.0±0.7	57.1±1.0	47.3±0.9	39.7±0.5	35.3±0.4	36.9±0.5
CLIP ₂	98.6±0.2	98.9±0.2	98.4±0.2	96.7±0.0	92.2±0.2	82.5±0.7	69.1±1.5	57.8±1.8	48.4±1.3	41.0±1.9	41.1±1.0
CLIP ₃	98.0±0.2	98.6±0.1	98.9±0.1	98.0±0.2	96.2±0.1	91.9±0.2	82.9±0.4	72.3±0.3	60.8±0.4	49.8±0.6	47.5±1.0
CLIP ₄	98.1±0.5	98.9±0.1	98.7±0.2	98.4±0.2	97.8±0.4	95.6±0.1	90.2±0.2	81.5±0.5	69.7±1.2	56.9±1.1	52.8±1.7

Table 23: ViT-B/32 results on ROTATEDMNIST.

	\mathcal{D}_0	\mathcal{D}_1	\mathcal{D}_2	\mathcal{D}_3	\mathcal{D}_4	\mathcal{D}_5	\mathcal{D}_6	\mathcal{D}_7	\mathcal{D}_8	\mathcal{D}_9	\mathcal{D}_{10}
IN ₀	97.7±0.1	92.7±0.2	85.5±0.5	75.5±1.1	61.7±1.3	47.8±1.0	37.1±1.0	29.4±1.5	23.9±0.5	20.2±0.8	17.2±1.1
IN ₁	98.1±0.5	96.2±0.6	88.2±0.8	78.2±0.8	65.9±0.9	52.4±0.5	39.7±1.1	32.1±1.2	25.9±0.6	20.1±0.6	18.1±1.0
IN ₂	97.1±0.6	96.6±0.3	90.7±0.4	81.7±0.8	71.7±0.3	59.4±1.6	47.8±1.3	37.5±1.2	29.0±0.9	23.0±1.9	19.3±0.7
IN ₃	97.4±0.5	96.0±0.7	92.6±1.4	85.5±0.7	75.5±1.5	65.6±1.5	52.8±0.6	41.7±0.4	34.1±0.8	26.4±1.6	20.8±1.0
IN ₄	97.7±0.2	96.0±0.3	93.1±0.8	85.6±2.4	76.6±1.8	67.1±0.4	56.6±0.9	47.1±2.0	39.2±0.7	31.7±0.6	25.2±0.8
CLIP ₀	93.8±0.3	87.7±0.3	77.4±0.9	62.5±1.5	46.5±1.4	34.5±0.4	26.5±1.6	19.9±0.9	16.0±1.4	14.5±1.1	12.4±1.0
CLIP ₁	93.4±0.9	90.5±0.7	81.4±1.0	67.8±0.2	50.0±1.2	33.5±1.9	23.2±0.7	15.7±0.4	12.2±0.4	10.2±0.6	8.7±0.5
CLIP ₂	93.9±0.7	90.9±0.6	84.4±0.5	76.2±1.1	64.2±1.1	49.3±1.9	34.2±1.6	23.9±1.2	18.0±1.0	13.2±0.8	10.9±1.0
CLIP ₃	93.0±0.8	90.5±1.2	85.5±1.8	77.6±1.3	68.6±1.0	56.9±1.1	43.8±0.5	31.3±1.7	22.7±1.2	17.2±0.9	14.0±0.4
CLIP ₄	92.9±1.3	91.0±1.6	84.8±1.5	76.8±1.8	71.0±1.4	60.6±1.8	52.3±1.3	40.7±1.1	32.3±1.7	24.4±0.7	19.0±0.6

Table 24: ResNet-50 results on NOISYIMAGENET15.

	\mathcal{D}_0	\mathcal{D}_1	\mathcal{D}_2	\mathcal{D}_3	\mathcal{D}_4	\mathcal{D}_5	\mathcal{D}_6	\mathcal{D}_7	\mathcal{D}_8	\mathcal{D}_9	\mathcal{D}_{10}
IN ₀	98.3±0.3	92.9±0.2	92.2±0.7	90.9±0.3	89.3±0.3	85.9±0.8	82.2±0.4	77.9±0.5	71.8±0.4	67.7±1.5	61.5±0.9
IN ₁	98.4±0.3	98.4±0.3	92.6±0.4	91.5±0.3	89.0±0.7	86.6±0.6	82.6±0.8	76.7±1.4	72.3±1.3	65.7±1.0	59.6±1.2
IN ₂	98.2±0.4	98.3±0.3	96.9±0.3	90.7±0.2	89.0±0.3	87.1±0.6	81.4±0.9	77.8±1.2	71.3±0.2	64.2±0.9	57.7±2.0
IN ₃	98.5±0.5	98.7±0.6	97.4±0.8	95.9±0.8	89.9±0.3	86.7±0.9	83.2±0.3	79.8±0.9	73.7±0.8	67.8±0.8	60.9±0.5
IN ₄	98.8±0.3	98.5±0.8	97.6±0.2	95.5±1.2	92.9±1.4	87.2±0.5	84.7±0.1	79.2±1.3	75.3±0.7	69.1±0.8	62.6±0.6
CLIP ₀	94.4±0.1	92.1±0.3	88.8±0.6	81.9±0.4	72.7±1.1	62.9±1.0	53.7±1.0	45.5±1.1	39.7±1.7	32.3±1.8	27.6±0.7
CLIP ₁	94.6±0.5	93.7±0.4	90.0±0.4	85.5±0.6	76.6±1.7	69.7±2.6	60.2±2.2	50.9±4.0	42.4±3.2	35.0±1.8	28.9±2.9
CLIP ₂	94.5±0.8	93.8±0.6	91.8±0.6	87.0±0.4	81.4±0.9	73.5±0.7	65.4±2.3	55.0±1.3	44.5±1.7	38.1±1.4	30.3±1.5
CLIP ₃	94.8±1.1	94.6±0.7	90.7±0.6	87.4±1.2	82.3±0.9	75.7±1.0	69.1±0.6	60.8±0.6	51.9±1.9	42.8±1.3	35.7±0.8
CLIP ₄	94.4±1.3	94.2±1.0	90.4±0.5	86.2±1.5	82.1±1.4	77.7±1.1	71.5±0.2	64.4±1.5	55.2±1.1	46.6±1.3	40.7±1.8

Table 25: ViT-B/32 results on NOISYIMAGENET15.

	\mathcal{D}_0	\mathcal{D}_1	\mathcal{D}_2	\mathcal{D}_3	\mathcal{D}_4	\mathcal{D}_5	\mathcal{D}_6	\mathcal{D}_7	\mathcal{D}_8	\mathcal{D}_9	\mathcal{D}_{10}
IN ₀	97.9±0.2	93.8±0.5	91.6±0.3	87.9±0.5	84.1±0.5	78.5±1.3	71.5±2.0	64.2±2.8	55.8±3.6	43.4±4.6	32.9±2.9
IN ₁	98.0±0.3	97.4±0.4	92.8±0.3	90.3±0.5	87.9±0.3	82.9±0.3	76.4±1.2	69.6±1.5	59.3±2.0	50.3±1.8	39.7±1.6
IN ₂	96.8±0.4	97.3±0.2	96.4±0.4	91.4±0.3	89.4±0.1	85.6±0.5	81.0±0.8	75.7±0.5	67.4±1.4	57.0±1.1	45.8±0.8
IN ₃	97.4±0.6	96.7±0.7	97.6±0.7	95.3±1.0	89.7±0.3	86.4±0.6	82.5±0.2	78.1±0.4	69.0±0.9	58.1±0.9	47.4±0.3
IN ₄	97.3±0.4	97.4±1.1	96.6±0.5	95.3±0.3	94.6±0.9	87.1±0.4	82.6±0.3	78.4±0.9	69.9±0.9	58.2±1.2	46.5±2.1
CLIP ₀	93.5±0.3	90.8±0.7	89.2±0.5	86.1±0.8	83.0±0.2	78.3±0.8	71.3±1.5	60.2±1.5	49.2±1.4	38.0±1.9	27.3±2.5
CLIP ₁	93.3±0.3	92.4±0.2	90.7±0.4	88.8±0.5	86.0±0.7	81.3±0.4	73.8±0.5	64.9±0.9	54.5±1.2	44.3±2.2	34.3±1.3
CLIP ₂	93.7±1.4	92.7±0.8	91.8±1.1	88.7±0.1	85.7±0.9	81.7±0.7	74.1±0.4	63.6±1.6	52.5±1.2	42.1±0.8	31.1±2.1
CLIP ₃	93.2±0.3	92.2±0.9	92.3±0.3	91.0±2.0	86.4±0.7	83.0±0.4	75.8±0.3	66.3±1.2	55.4±1.9	44.7±1.0	33.1±1.8
CLIP ₄	92.2±0.5	93.4±0.9	91.5±0.2	91.3±1.4	89.4±1.0	83.5±0.2	77.3±0.9	68.4±0.3	58.0±0.6	46.7±0.5	34.2±1.1

Table 26: ResNet-50 results on LR-IMAGENET15.

	\mathcal{D}_0	\mathcal{D}_1	\mathcal{D}_2	\mathcal{D}_3	\mathcal{D}_4	\mathcal{D}_5	\mathcal{D}_6	\mathcal{D}_7	\mathcal{D}_8	\mathcal{D}_9	\mathcal{D}_{10}
IN ₀	98.5±0.2	94.2±0.1	93.6±0.2	92.2±0.1	90.3±0.2	88.4±0.3	85.6±0.1	81.2±0.8	77.0±1.1	70.9±1.3	59.8±1.4
IN ₁	98.1±0.4	98.6±0.3	93.7±0.1	92.4±0.1	91.1±0.3	89.2±0.2	86.5±0.2	83.4±0.3	79.5±0.2	73.6±0.3	62.7±0.5
IN ₂	98.2±0.8	98.3±0.4	97.9±0.3	92.8±0.2	91.0±0.2	89.2±0.2	86.9±0.2	84.3±0.2	80.5±0.0	74.6±0.3	63.6±0.7
IN ₃	98.5±1.3	98.3±0.8	98.4±0.0	98.2±0.2	91.3±0.4	90.0±0.5	87.4±0.3	83.9±0.0	80.4±0.4	74.4±0.3	63.6±0.5
IN ₄	98.9±0.3	98.3±0.5	98.3±0.7	98.3±0.1	97.6±0.6	90.6±0.1	88.1±0.3	85.1±0.4	81.1±0.5	75.2±0.3	66.2±0.7
CLIP ₀	94.5±0.3	93.3±0.1	92.1±0.4	91.0±0.1	89.9±0.4	88.3±0.3	85.7±0.3	80.0±0.7	73.6±1.0	62.7±1.2	54.1±1.9
CLIP ₁	94.1±0.7	94.8±0.5	92.4±0.1	91.6±0.1	90.3±0.1	88.6±0.8	86.2±1.8	79.3±1.6	73.0±2.1	63.2±1.4	54.0±0.8
CLIP ₂	94.8±1.1	94.2±1.0	93.9±1.0	92.1±0.1	91.0±0.2	89.9±0.3	86.6±0.3	79.8±0.3	74.0±0.3	64.6±0.5	55.2±0.5
CLIP ₃	94.0±0.7	94.4±0.9	94.9±0.3	94.1±0.8	90.8±0.1	89.3±0.4	86.0±1.2	80.0±0.9	73.7±1.4	64.4±1.1	55.1±2.6
CLIP ₄	94.9±0.4	93.6±1.0	94.3±1.2	92.9±0.2	92.1±0.6	90.0±0.2	86.4±0.3	81.0±0.7	75.6±0.6	64.6±0.3	56.2±0.3

Table 27: ViT-B/32 results on LR-IMAGENET15.