# Covariate Shift Detection via Domain Interpolation Sensitivity

**Tejas Gokhale\*** Arizona State University tgokhale@asu.edu

Joshua Feinglass\* Arizona State University joshua.feinglass@asu.edu Yezhou Yang Arizona State University yz.yang@asu.edu

#### **Abstract**

Covariate shift is a major roadblock in the reliability of image classifiers in the real world. Work on covariate shift has been focused on training classifiers to adapt or generalize to unseen domains. However, for transparent decision making, it is equally desirable to develop *covariate shift detection* methods that can indicate whether or not a test image belongs to an unseen domain. In this paper, we introduce a benchmark for covariate shift detection (CSD), that builds upon and complements previous work on domain generalization. We use state-of-the-art OOD detection<sup>1</sup> methods as baselines and find them to be worse than simple confidence-based methods on our CSD benchmark. We propose an interpolation-based technique, Domain Interpolation Sensitivity (DIS), based on the simple hypothesis that interpolation between the test input and randomly sampled inputs from the training domain, offers sufficient information to distinguish between the training domain and unseen domains under covariate shift. DIS surpasses all OOD detection baselines for CSD on multiple domain generalization benchmarks.

## 1 Introduction

Machine learning models such as image classifiers are being increasingly deployed in real-world settings. Covariate shift is a commonly occurring phenomena, where test images are from the same categories as the training data, but undergo a shift in terms of style. For instance, the training data may contain images taken during the day under sunny conditions, but the classifier may encounter nighttime images or foggy or rainy images. Models trained under the empirical risk minimization [36] paradigm can only offer performance guarantees under the *i.i.d.* setting, and are known to fail under various types of covariate shift [33, 17, 2].

To mitigate the risks associated with covariate shift, domain generalization algorithms have been developed [1, 38, 40, 10]. However, improving the accuracy of classifiers on unseen domains cannot be the only criteria for reliable decision making – for transparency, methods that *detect* covariate shift should also be investigated. Unfortunately, this aspect of reliable decision making has not been previously explored. In this paper, we investigate covariate shift detection (CSD) for image classifiers.

Methods for detecting "out-of-distribution" (OOD) test examples have been previously developed [18, 24, 23, 20, 3]. However it is important to note that OOD detection algorithms are designed to detect *novel categories* at test time. In this paper, we are interested in detecting covariate shift (i.e. detecting test inputs that belong to a previously unseen domain, but one of the classes that the classifier is trained on). Towards this end, we develop a new benchmark for covariate shift detection. We utilize common domain generalization benchmarks and train the classifier on one of the domains, and benchmark CSD methods' performance in detecting images that belong to other domains. An example is shown

<sup>\*</sup>These authors contributed equally to this work.

<sup>&</sup>lt;sup>1</sup>Recent literature uses the term "OOD detection" to refer to novel category detection.

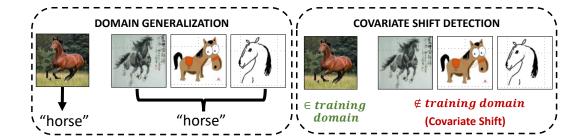


Figure 1: Domain generalization (DG) and covariate shift detection (CSD) are both important, but orthogonal aspects of robustness evaluation. While the aim of DG is to predict the correct label for inputs from unseen domains, the aim of CSD is to detect unseen domains – i.e. detect covariate shift in test inputs.

in Figure 1 – we compare the two dimensions of reliable predictions under covariate shift. Domain generalization algorithms are trained with the aim of making accurate predictions for seen as well as unseen domains. However in cases where accuracy under domain shift is low, it is equally important to detect, or flag, cases where there may be a covariate shift for safe and reliable use of classifiers.

Interpolation between training examples has been shown to provide unique information suitable for model regularization [31, 21]. Motivated by this, we hypothesize that interpolating test inputs with training inputs offers sufficient information to distinguish between in-domain and out-of-domain examples. We use test-time interpolations as a powerful tool for understanding the behavior of pre-trained classifiers and detecting covariate shift. Our motivation is as follows: if interpolations during training benefit model generalization, then interpolating test inputs with training inputs could help us understand how the model might perform on unseen test distributions.

We develop a method, named "Domain Interpolation Sensitivity" (DIS), that achieves state-of-the-art results on covariate shift detection for image classification and text classification. We find that methods developed for OOD detection (novel category detection) underperform on the CSD benchmark. Suprisingly we find that methods that perform better on OOD detection benchmarks than the maximum softmax probability (MSP) baseline by Hendrycks et al. [18], perform much worse than MSP on our CSD benchmarks.

Our contributions are summarized below:

- We study covariate shift detection (CSD) as a mechanism for improving the reliability of classifier predictions; CSD is designed to complement domain generalization as a robustness metric.
- We propose CSD benchmarks that are derived from existing benchmarks for domain generalization.
- We develop a interpolation-based technique that outperforms existing outlier and OOD detection methods on four CSD benchmarks (three for image classification, and one for text classification).

# 2 Covariate Shift Detection

We will consider classification tasks, for which a neural network f is trained on a dataset  $\mathcal{D}_{in}$  containing labeled input–output pairs  $(\mathbf{x},\mathbf{y})$ , with inputs  $\mathbf{x} \in \mathcal{X}_{in}$  and outputs  $\mathbf{y} \in \mathcal{Y}_{in}$ . Let  $\mathcal{D}_{out}$  denote previously unseen data. The nature of the shift between  $\mathcal{D}_{in}$  and  $\mathcal{D}_{out}$  can take multiple forms. One such type of distribution shift is the presence of novel categories in  $\mathcal{D}_{out}$ , i.e. if  $(\mathbf{x},\mathbf{y}) \in \mathcal{D}_{out}$ , then the categorical label of  $\mathbf{x},\mathbf{y} \notin \mathcal{Y}_{out}$ . An example of this phenomena of novel categories is if  $\mathcal{D}_{in}$  is a dataset for cat-dog classification, whereas  $\mathcal{D}_{out}$  contains images of handwritten digits. Another type of distribution shift can occur when the categorical space remains the same, but domain  $\mathcal{X}_{out}$  undergoes a covariate shift, i.e.  $p_{in}(\mathbf{x}) \neq p_{out}(\mathbf{x})$ . For example, a covariate shift exists between  $\mathcal{D}_{in}$  and  $\mathcal{D}_{out}$  if  $\mathcal{D}_{in}$  is a set of real cat-dog images while  $\mathcal{D}_{out}$  contains cartoons or sketches of cats and dogs. In this paper, we will consider covariate shift.

Scoring Function for Covariate Shift Detection. Covariate shift detection can be formulated as a binary classification task. Given a classifier f trained on distribution  $\mathcal{D}_{in}$ , the goal is to design a

estimator g that estimates whether or not a test input lies within the training domain.

$$g(x) = \begin{cases} 1 & \text{if } S(x) \ge \gamma \quad \text{(in-domain)} \\ 0 & \text{if } S(x) < \gamma \quad \text{(covariate shift)} \end{cases}$$
 (1)

The threshold  $\gamma$  is chosen such that 95% of in-domain data is correctly classified by Eq. 1 The choice of the scoring function S is the key to improving covariate shift detection. Previous approaches for OOD detection have utilized the model's outputs (for eg. the maximum softmax probability [18], or energy of the softmax output [24]), or model's gradient space [20]. In this work, we develop a scoring function S(x) by leveraging the interpolation of the test input with training inputs.

Benchmarking Covariate Shift Detection. We leverage existing domain generalization datasets for benchmarking CSD. Specifically, we operate under the single-source domain generalization setting [38], where the classifier is trained only on one domain and tested on all domains within the dataset. This is illustrated in Figure 1 which shows domain generalization on the PACS [22] dataset – the classifier is trained on real-world photos (source domain) and tested on all domains including photos, art-paintings, cartoons, and sketches. Given a classifier trained on the source domain  $\mathcal{D}_i n$ , the goal of a covariate shift detection algorithm is to use a scoring function S(x) to estimate whether or not x belongs to the source domain or not. Thus, for any domain generalization dataset, we can compare performance of CSD algorithms on the corresponding CSD benchmark.

# 3 Domain Interpolation Sensitivity

In this section, we describe our method for covariate shift detection using test-time input interpolation.

**Test-Time Input Interpolation.** Consider a randomly sampled training input  $x_i \in \mathcal{X}_{in}$  and a test input x. We define the interpolation of  $x_i$  and x to be  $\hat{x} = h(x_i, x, \epsilon)$  for a mixing coefficient  $\epsilon \in [0, 1]$ . Note that  $h(x_i, x, 0) = x$  and  $h(x_i, x, 1) = x_i$ . In practice, h can be implemented in multiple ways; for image classification tasks, we use a simple pixel-wise convex combination. For text classification, we use a token-wise swapping.

$$h_{nixelwise}(\mathbf{x}_i, \mathbf{x}, \epsilon) = \epsilon \mathbf{x}_i + (1 - \epsilon)\mathbf{x}$$
 (2)

**Domain Interpolation Sensitivity.** Let  $[0,\epsilon,2\epsilon\dots T\epsilon]$  be an increasing sequence of mixing coefficients such that  $0\leq\epsilon\leq T\epsilon\leq 1$ . We generate a sequence of interpolated images

$$X_i = [h(x_i, x, 0), h(x_i, x, \epsilon), \dots h(x_i, x, T\epsilon)].$$
(3)

We obtain a corresponding sequence of softmax predictions probabilities from model f as:

$$Y_i = f(X_i) = [f(h(x_i, x, 0)), f(h(x_i, x, \epsilon_1)), \dots f(h(x_i, x, \epsilon_T))].$$
 (4)

Note that we can generate  $Y_i$  for each choice of training exemplar  $X_i$ . By using n training exemplars, we can obtain an average prediction sequence  $\bar{Y}$ :

$$\bar{Y} = \left[\frac{1}{n}\sum_{i=1}^{n} f(h(x_i, x, 0)), \frac{1}{n}\sum_{i=1}^{n} f(h(x_i, x, \epsilon_1)), \dots \frac{1}{n}\sum_{i=1}^{n} f(h(x_i, x, \epsilon_T))\right].$$
 (5)

Let  $c = \underset{\mathcal{Y}_{in}}{argmax} f(\mathbf{x})$  be the predicted category for the test input. Then, the domain interpolation sensitivity curve is defined as the softmax probability of c in each element of  $\bar{\mathbf{Y}}$ .

Once we've obtained the above DIS curves, we use the area under the DIS curve as the scoring function S(x) for covariate shift detection, i.e.  $S(x) = AUC(\bar{Y})$ .

# 4 Experiments

**Baselines.** We use widely adopted OOD detection approaches MSP [18], Energy [24], ODIN [23], and GradNorm [20] as our baselines. These methods do not rely on any additional training or other modifications to the classifier used for detection and are thus comparable to our own approach.

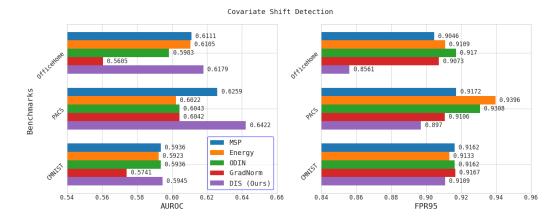


Figure 2: Summary of results on covariate shift detection benchmarks for image classification, in terms of AUROC (higher value is better) and FPR95 (lower value is better). Detailed results can be found in the appendix. The key observation is that recent OOD methods perform worse than the baseline MSP [18], in terms of both AUROC and FPR95. Our DIS method improves CSD detection performance on all three benchmarks.

Method	<b>Rotten Tomatoes</b>	IMDB	Amazon Reviews	Average
MSP [18]	0.7966 / 0.7337	0.6718 / 0.8460	0.6133 / 0.889	0.6939 / 0.8229
Energy [24]	0.7827 / 0.7349	0.6701 / 0.8390	0.6037 / 0.881	0.6855 / 0.8183
GradNorm [20]	0.7894 / 0.7349	0.6692 / 0.8480	0.6086 / 0.8930	0.6891 / 0.8253
DIS (Ours)	0.9209 / 0.5917	0.6224 / 0.8960	0.7184 / 0.8690	0.7539 / 0.7856

Table 1: Covariate shift detection performance on the proposed text classification benchmark. Results are shown as AUROC ↑ / FPR95 ↓. We observe that our method is highly consistent across both tested modalities.

**Metrics.** For comparative evaluation of our method against the baselines, we use two standard metrics [18]; these are: (i) **AUROC**: area under the ROC curve [5], (ii) **FPR95**: false positive rate on the OOD set when the true positive rate on the ID set is 95%.

**Datasets.** We consider three image classification benchmarks: PACS [22], OfficeHome [37], and ColoredMNIST [1], as well as a proposed text review classification benchmark inspired by [19]. PACS contains four domains: photos, art-paintings, cartoons, and sketches. OfficeHome contains four domains: real images, art, clipart, and product images. ColoredMNIST contains three domains of digit images with varying degrees of spurious correlations (+90, +80, -90) between the digit and color. The proposed text classification benchmark contains four sentiment analysis domains collected from different review websites: Yelp, Rotten Tomatoes, IMDB [25], and Amazon [27, 15].

For the image classification benchmarks, we follow the training protocol from DomainBed [13] and train a ResNet-18 model [14] on 'Photos' for PACS, 'Real' for OfficeHome and the domain with "+90" spurious correlation for ColoredMNIST. For our text benchmark, we use BERT [6] with a classification head fine-tuned on the Yelp domain and 1000 randomly selected examples from the test sets of the other domains for covariate shift detection.

**Hyperparameters.** For image classification experiments, we use n=16 exemplars, step size  $\epsilon=0.05$ , and number of interpolation steps T=4. For text classification, we use n=16 exemplars, step size  $\epsilon=1$  token, and number of interpolation steps T equal to a quarter of the total token length.

**Results.** Our results on the image classification and text classification benchmarks reveal the strength of the interpolation-based method, with consistent improvements both in terms of AUROC and FPR95 on both tasks. Figure 2 summarizes our results on the image classification CSD benchmarks. Interestingly, we observe that Energy, ODIN, and GradNorm, which have been shown to be better than MSP in the OOD detection literature, are in fact, much worse than MSP on all three benchmarks (OfficeHome, PACS, ColoredMNIST). DIS outperforms all baselines. Table 1 shows results on the text classification benchmark. The observations are similar – sophisticated OOD detection methods perform worse than the baseline MSP, while DIS outperforms all methods.

# 5 Related Work

**Distributional Robustness.** Several dimensions of distribution robustness have been studied, which can be broadly classified into adversarial robustness [12, 26], natural distributional robustness, and spurious correlations. Work on natural distributional robustness includes conditions such as common corruptions [17], variations along style [16], geometric [39] and attribute-level shift [9], different dataset sources [37, 22]. Spurious correlations of features such as background [2, 32] or texture [8] with label space have been studied. In terms of label shift, anomaly/outlier detection, novelty detection, open-set recognition, and OOD detection have been studied [41].

Correlation between ID and OOD performance. It is well known that models tested on data with covariate shift suffer a drop in performance compared to in-domain (ID) accuracy. Recently, there have been several studies that find a positive correlation between ID and OOD performance for tasks in both computer vision [29] and natural language processing [28]. However, Teney et al.[34] show that under certain real-world conditions, a negative correlation might exist, i.e. a decrease in ID accuracy may benefit OOD performance. Moayeri et al.[30] show that there are trade-offs between adversarial and natural distributional robustness. There is also empirical evidence [11] that suggests that data modification techniques (for instance, data augmentation or data filtering [4]) may have a negative impact on adversarial robustness. With the context of these findings, there is a large gap in our understanding of different robustness settings – characterizing distribution shift in different ways is therefore crucial as a model selection criteria. Our work aims to aid investigations in this direction.

**Interpolation for Model Selection.** Interpolations have been extensively used for representation learning [43, 42, 31] for training robust classifiers. Bhattacharjee et al.[3] used interpolation of inputs during training for novel category detection. SMURF [7] used token interpolation between input texts and a noise process to estimate the robustness of language models based on the average monotonicity of a perplexity measure.

## 6 Outlook

In this work we introduced covariate shift detection benchmarks to study a complementary measure of model robustness. Our experiments reveal that for the CSD task, sophisticated OOD detection methods are worse than even the simple MSP baseline. We present a simple interpolation-based detection technique that surpasses all baselines on multiple CSD benchmarks on both image classification and text classification tasks. The results are promising and suggest that interpolation between training and test inputs can be a powerful tool for understanding and interpreting classification decisions as well as detecting outliers and covariate shift. We believe that test-time interpolation could also be useful for uncertainty quantification – recent results [35] show how anchoring (a variant of interpolation) can be used during training for this purpose. In the future, we expect theoretical insights to emerge to complement our empirical findings with DIS.

## Acknowledgements

This work was supported by NSF CPS grant #2038666 and RI grant #2132724, and Amazon AWS Machine Learning Research Award.

#### References

- [1] Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization. *arXiv* preprint arXiv:1907.02893, 2019.
- [2] Sara Beery, Grant Van Horn, and Pietro Perona. Recognition in terra incognita. In *Proceedings of the European conference on computer vision (ECCV)*, pages 456–473, 2018.
- [3] Supritam Bhattacharjee, Devraj Mandal, and Soma Biswas. Multi-class novelty detection using mix-up technique. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, March 2020.
- [4] Ronan Le Bras, Swabha Swayamdipta, Chandra Bhagavatula, Rowan Zellers, Matthew E. Peters, Ashish Sabharwal, and Yejin Choi. Adversarial filters of dataset biases. In *Proceedings of the 37th International*

- Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event, volume 119 of Proceedings of Machine Learning Research, pages 1078–1088. PMLR, 2020.
- [5] Jesse Davis and Mark Goadrich. The relationship between precision-recall and roc curves. In *Proceedings* of the 23rd international conference on Machine learning, pages 233–240, 2006.
- [6] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4171–4186, Minneapolis, Minnesota, 2019. Association for Computational Linguistics.
- [7] Joshua Feinglass and Yezhou Yang. SMURF: SeMantic and linguistic UndeRstanding fusion for caption evaluation via typicality analysis. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 2250–2260, Online, August 2021. Association for Computational Linguistics.
- [8] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. In *International Conference on Learning Representations*, 2018.
- [9] Tejas Gokhale, Rushil Anirudh, Bhavya Kailkhura, Jayaraman J Thiagarajan, Chitta Baral, and Yezhou Yang. Attribute-guided adversarial training for robustness to natural perturbations. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 7574–7582, 2021.
- [10] Tejas Gokhale, Rushil Anirudh, Jayaraman J Thiagarajan, Bhavya Kailkhura, Chitta Baral, and Yezhou Yang. Improving diversity with adversarially learned transformations for domain generalization. In IEEE/CVF Winter Conference on Applications of Computer Vision, 2023.
- [11] Tejas Gokhale, Swaroop Mishra, Man Luo, Bhavdeep Sachdeva, and Chitta Baral. Generalized but not Robust? comparing the effects of data modification methods on out-of-domain generalization and adversarial robustness. In Findings of the Association for Computational Linguistics: ACL 2022, pages 2705–2718. Association for Computational Linguistics, May 2022.
- [12] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In Yoshua Bengio and Yann LeCun, editors, 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings, 2015.
- [13] Ishaan Gulrajani and David Lopez-Paz. In search of lost domain generalization. In *International Conference on Learning Representations*, 2021.
- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In 2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016, pages 770–778. IEEE Computer Society, 2016.
- [15] Ruining He and Julian McAuley. Ups and downs: Modeling the visual evolution of fashion trends with one-class collaborative filtering. In *Proceedings of the 25th International Conference on World Wide Web*, WWW '16, page 507–517, Republic and Canton of Geneva, CHE, 2016. International World Wide Web Conferences Steering Committee.
- [16] 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*, pages 8340–8349, 2021.
- [17] Dan Hendrycks and Thomas G. Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net, 2019.
- [18] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *ArXiv*, abs/1610.02136, 2017.
- [19] Dan Hendrycks, Xiaoyuan Liu, Eric Wallace, Adam Dziedzic, Rishabh Krishnan, and Dawn Song. Pretrained transformers improve out-of-distribution robustness. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2744–2751, Online, 2020. Association for Computational Linguistics.
- [20] Rui Huang, Andrew Geng, and Yixuan Li. On the importance of gradients for detecting distributional shifts in the wild. In Advances in Neural Information Processing Systems, 2021.

- [21] David Krueger, Ethan Caballero, Joern-Henrik Jacobsen, Amy Zhang, Jonathan Binas, Dinghuai Zhang, Remi Le Priol, and Aaron Courville. Out-of-distribution generalization via risk extrapolation (rex). In International Conference on Machine Learning, pages 5815–5826. PMLR, 2021.
- [22] Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M. Hospedales. Deeper, broader and artier domain generalization. In *IEEE International Conference on Computer Vision*, *ICCV 2017*, *Venice*, *Italy*, *October 22-29*, 2017, pages 5543–5551. IEEE Computer Society, 2017.
- [23] Shiyu Liang, Yixuan Li, and R Srikant. Enhancing the reliability of out-of-distribution image detection in neural networks. *International Conference on Learning Representations (ICLR)*, 2018.
- [24] Weitang Liu, Xiaoyun Wang, John Owens, and Yixuan Li. Energy-based out-of-distribution detection. Advances in Neural Information Processing Systems (NeurIPS), 2020.
- [25] Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. Learning word vectors for sentiment analysis. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 142–150, Portland, Oregon, USA, June 2011. Association for Computational Linguistics.
- [26] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings. OpenReview.net, 2018.
- [27] Julian McAuley, Christopher Targett, Qinfeng Shi, and Anton van den Hengel. Image-based recommendations on styles and substitutes. In *Proceedings of the 38th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '15, page 43–52, New York, NY, USA, 2015. Association for Computing Machinery.
- [28] John Miller, Karl Krauth, Benjamin Recht, and Ludwig Schmidt. The effect of natural distribution shift on question answering models. In *International Conference on Machine Learning*, pages 6905–6916. PMLR, 2020.
- [29] 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 *International Conference on Machine Learning*, pages 7721–7735. PMLR, 2021.
- [30] Mazda Moayeri, Kiarash Banihashem, and Soheil Feizi. Explicit tradeoffs between adversarial and natural distributional robustness. arXiv preprint arXiv:2209.07592, 2022.
- [31] 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*, pages 8690–8699, 2021.
- [32] Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. In *International Conference on Learning Representations*, 2020.
- [33] Rohan Taori, Achal Dave, Vaishaal Shankar, Nicholas Carlini, Benjamin Recht, and Ludwig Schmidt. Measuring robustness to natural distribution shifts in image classification. In Advances in Neural Information Processing Systems, volume 33, pages 18583–18599, 2020.
- [34] Damien Teney, Yong Lin, Seong Joon Oh, and Ehsan Abbasnejad. Id and ood performance are sometimes inversely correlated on real-world datasets. *arXiv preprint arXiv:2209.00613*, 2022.
- [35] Jayaraman J Thiagarajan, Rushil Anirudh, Vivek Narayanaswamy, and Peer-Timo Bremer. Single model uncertainty estimation via stochastic data centering. *NeurIPS*, 2022.
- [36] Vladimir N Vapnik and A Chervonenkis. The necessary and sufficient conditions for consistency of the method of empirical risk minimization. *Pattern Recognition and Image Analysis*, 1(3):284–305, 1991.
- [37] Hemanth Venkateswara, Jose Eusebio, Shayok Chakraborty, and Sethuraman Panchanathan. Deep hashing network for unsupervised domain adaptation. In 2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017, pages 5385–5394. IEEE Computer Society, 2017.

- [38] Riccardo Volpi, Hongseok Namkoong, Ozan Sener, John C. Duchi, Vittorio Murino, and Silvio Savarese. Generalizing to unseen domains via adversarial data augmentation. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada, pages 5339–5349, 2018.
- [39] Eric Wong and J Zico Kolter. Learning perturbation sets for robust machine learning. In *International Conference on Learning Representations*, 2020.
- [40] Zhenlin Xu, Deyi Liu, Junlin Yang, Colin Raffel, and Marc Niethammer. Robust and generalizable visual representation learning via random convolutions. In *International Conference on Learning Representations*, 2020.
- [41] Jingkang Yang, Kaiyang Zhou, Yixuan Li, and Ziwei Liu. Generalized out-of-distribution detection: A survey. arXiv preprint arXiv:2110.11334, 2021.
- [42] Sangdoo Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.
- [43] Hongyi Zhang, Moustapha Cissé, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 May 3, 2018, Conference Track Proceedings. OpenReview.net, 2018.

# **Appendix**

The tables below show detailed results on each image classification benchmark, split by domain. These correspond to the summarized results in Figure 2.

Method	Art	Clipart	Product	Average
MSP [18]	0.6863 / 0.8763	0.6379 / 0.9200	0.5534 / 0.9552	0.6259 / 0.9172
Energy [24]	0.6716 / 0.9084	0.6016 / 0.9448	0.5334 / 0.9656	0.6022 / 0.9396
ODIN [23]	0.6753 / 0.8946	0.6074 / 0.9414	0.5301 / 0.9564	0.6043 / 0.9308
GN [20]	0.6479 / 0.8694	0.6106 / 0.9256	0.5542 / 0.9369	0.6042 / 0.9106
Ours	0.7093 / 0.8339	0.6529 / 0.9087	0.5645 / 0.9483	0.6422 / 0.8970

Table 2: Covariate Shift Detection performance on the OfficeHome benchmark. All methods use the same ResNet classifier trained on the "Real" domain. Results are shown as AUROC ↑ / FPR95 ↓.

Method	Art	Cartoon	Sketch	Average
MSP [18]	0.6652 / 0.8868	0.4159 / 0.9760	0.7522 / 0.8510	0.6111 / 0.9046
Energy [24]	0.6708 / 0.8782	0.4252 / 0.9641	0.7356 / 0.8904	0.6105 / 0.9109
ODIN [23]	0.6571 / 0.9060	0.4098 / 0.9431	0.7279 / 0.9019	0.5983 / 0.9170
GN [20]	0.6298 / 0.8782	0.3980 / 0.9521	0.6536 / 0.8917	0.5605 / 0.9073
Ours	0.6693 / 0.8227	0.4278 / 0.9341	0.7567 / 0.8115	0.6179 / 0.8561

Table 3: Covariate Shift Detection performance on the PACS benchmark. All methods use the same ResNet classifier trained on the "Photos" domain. Results are shown as AUROC  $\uparrow$  / FPR95  $\downarrow$ 

Method	+80	-90	Average
MSP [18]	0.5225 / 0.9378	0.6647 / 0.8946	0.5936 / 0.9162
Energy [24]	0.5225 / 0.9361	0.6620 / 0.8905	0.5923 / 0.9133
ODIN [23]	0.5225 / 0.9378	0.6647 / 0.8946	0.5936 / 0.9162
GN [20]	0.5213 / 0.9381	0.6269 / 0.8954	0.5741 / 0.9167
Ours	0.5242 / 0.9379	0.6649 / 0.8841	0.5945 / 0.9109

Table 4: Covariate Shift Detection performance on the ColoredMNIST benchmark. All methods use the same CNN classifier and results are shown as AUROC  $\uparrow$  / FPR95  $\downarrow$