
Diversity Boosted Learning for Domain Generalization with A Large Number of Domains

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Machine learning algorithms minimizing the average training loss typically suf-
2 fer from poor generalization performance. It inspires various works for domain
3 generalization (DG), among which a series of methods work by $O(n^2)$ pairwise
4 domain operations with n domains, where each one is often costly. Moreover,
5 while a common objective in the DG literature is to learn invariant representations
6 against spurious correlations induced by domains, we point out its insufficiency
7 and highlight the importance of alleviating spurious correlations caused by *objects*.
8 Based on the observation that diversity helps mitigate spurious correlations, we
9 propose a Diversity boosted twO-level saMplIng framework (DOMI) to efficiently
10 sample the most informative ones among a large number of domains and data
11 points. We show that DOMI helps train robust models against spurious correlations
12 from both domain-side and object-side, substantially enhancing the performance
13 of five backbone DG algorithms on Rotated MNIST and Rotated Fashion MNIST.

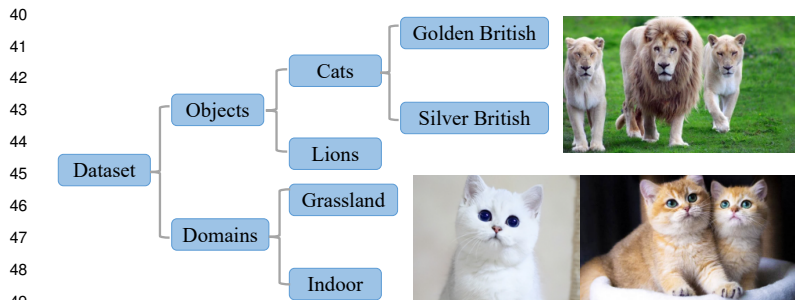
14 1 Introduction

15 The effectiveness of machine learning algorithms that minimize the average training loss relies on the
16 IID hypothesis. However, distributional shifts between test and training data are usually inevitable.
17 Under such circumstances, models trained by minimizing the average training loss are prone to sink
18 into spurious correlations. These misleading heuristics only work well on some data distributions
19 but can not be generalized to others that may appear in the test set. In domain generalization (DG)
20 tasks, the data distributions are denoted as different domains. The goal is to learn a model that can
21 generalize well to unseen ones after training on several domains. While lots of methods have been
22 derived to efficiently achieve this goal and show good performance, there are two main drawbacks.

23 **Scalability.** With an unprecedented amount of applicable data nowadays, many datasets contain a
24 tremendous amount of domains, massive data in each domain, or both. For instance, DrugOOD (Ji
25 et al., 2022) is an out-of-distribution dataset curator and benchmark for AI-aided drug discovery.
26 Datasets of DrugOOD contain hundreds to tens of thousands of domains. In addition to raw data with
27 multitudinous domains, domain augmentation, leveraged to improve the robustness of models in DG
28 tasks, can also lead to a significant increase in the number of domains. For example, HRM (Liu et al.,
29 2021) generates heterogeneous domains to help exclude variant features, favoring invariant learning.
30 Under such circumstances, training on the whole dataset in each epoch is computationally prohibitive,
31 especially for methods training by pairwise domain operations, of which the computational complexity
32 is $O(n^2)$ with n training domains.

33 **Objective.** Numerous works in the DG field focus entirely on excluding or alleviating domain-side
34 impacts. A general assumption in the DG field is that data in different domains share some “stable”
35 features to form causal correlations. And a large branch of studies holds that the relationship between

36 these “stable” features and the outputs is **domain-independent** given certain conditions (Long et al.,
 37 2015; Hoffman et al., 2018; Zhao et al., 2018, 2019; Mahajan et al., 2021). We state that this objective
 38 is insufficient, and a simple counterexample is given as follows. We highlight the importance of
 39 mitigating spurious correlations caused by the objects for training a robust model.



40
41
42
43
44
45
46
47
48
49
50 Figure 1: The training set of the counterexample. Cats are mainly silver British
 51 shorthair (body color of which is silvery white), rarely golden British shorthair
 52 (tan), and lions are all tan. As for the background, most lions are on the grassland
 53 while most cats are indoors.

Suppose our learning task is training a model to distinguish between cats and lions. The composition of the training set is shown in Figure 1, and the domain here refers to the images’ backgrounds. In this example, the correlation between features corresponding to the body color of the objects and class labels is undoubtedly independent of domains.

54 Moreover, it helps get high accuracy in the training set by simply taking the tan objects as li-
 55 ons and the white ones as cats. Unfortunately, if this correlation is mistaken for the causal correlation,
 56 the model is prone to poor performance once cat breed distribution shifts in the test set.

57 To tackle these two issues, we propose a diversity boosted two-level sampling framework named DOMI
 58 with the following major contributions: 1) To our best knowledge, this is the first paper to take impacts
 59 from the object side into account for achieving the goal of DG. 2) We propose DOMI, a diversity-
 60 boosted two-level sampling framework to select the most informative domains and data points for
 61 mitigating both domain-side and object-side impacts. 3) We demonstrate that DOMI substantially
 62 enhances the test accuracy of the backbone DG algorithms on two benchmarks.

63 2 Methods

64 We introduce our method DOMI by firstly presenting two key observations.

65 **Observation 1.** *Diverse domains of data help exclude spurious correlations.*

66 Consider a dataset $D_n = \{D^1, D^2, \dots, D^n\}$ which is a mixture of data $D^d = \{(x_i^d, y_i^d)\}_{i=1}^{n_d}$ where d
 67 is one domain of the ground set D ($|D| = n$), x_i^d and y_i^d are the i^{th} data and label from domain d
 68 respectively, and n_d is the number of data points in D_d . Suppose we now have dataset D_k consisting
 69 of k domains. On D_k , the distribution of data is $P^k(X, Y)$. A “good” set denoted by C_k is a set
 70 containing “good” correlations that get high accuracy on D_k . The set of causal correlations is C .
 71 $C \subseteq C_k$ since causal correlations can definitely get good performance but “good” correlations for the
 72 k domains may not be held in other domains, i.e., spurious correlations. The goal is to exclude as
 73 many spurious correlations as possible.

74 Given another domain d_{k+1} to form dataset D_{k+1} together with the former k domains. The corre-
 75 sponding data distribution and the “good” set are $P^{k+1}(X, Y)$ and C_{k+1} , respectively. If $P^{k+1}(X, Y)$
 76 is close to $P^k(X, Y)$, then most of the correlations in C_k will still be “good” for D_{k+1} and thus
 77 preserved in C_{k+1} . Nevertheless, if d_{k+1} is a heterogeneous domain that can significantly change the
 78 distribution of data, then the “good” set after being constrained would be obviously smaller than the
 79 original one, i.e., $|C_{k+1}| \ll |C_k|$, showing that diverse domains help exclude spurious correlations
 80 and training on which helps obtain robust models.

81 We formally derive Proposition 1 to support Observation 1 that diversity helps mitigate spurious
 82 correlations, based on which DOMI is a diversity boosted sampling framework and the sampling
 83 scheme to obtain a heterogeneous subset is a critical part of DOMI. Determinantal Point Process (DPP)
 84 (Kulesza et al., 2012) sampling is a powerful diversity sampling method. Based on the similarity
 85 matrix (DPP kernel) among the samples, a draw from a DPP yields diversified subsets. Thus we
 86 incorporate DPP sampling into DOMI. As one option for the diversity sampling method in DOMI, DPP

87 sampling can also be substituted with other sampling methods, which will be left as an interesting
88 future direction.

89 2.1 invDANN

90 Domain-Adversarial Neural Networks (DANN) proposed by Ganin et al. (2016) is composed by
91 Featurizer, Classifier, and Discriminator. Featurizer extracts features of data samples, Classifier
92 learns to classify class labels of data, and Discriminator learns to discriminate domains. DANN set a
93 gradient reversal layer between Featurizer and Discriminator to ensure Featurizer captures object-side
94 features. Using the architecture of DANN, we let Classifier learn to classify domain labels while
95 Discriminator learns to discriminate class labels. As an inverse version of DANN, invDANN trains a
96 model whose Featurizer extracts only domain-side features, which serves as an important component
97 of the proposed method.

98 Observation 2. Excluding domain-induced spurious correlations is insufficient for achieving OOD
99 generalization under the setting of DG.

100 Figure 2 shows a structural causal model (SCM) that describes the data-generating process for the
101 domain generalization task with object-side spurious correlations.

102 The SCM divides data into two parts: domain-side and object-side. \bar{x} of domain-side is the reason for domain-induced spurious
103 correlations. For the object-side, the feature is further divided into \underline{x} (causal features) and \tilde{x} where \tilde{x} is the reason for object-induced
104 spurious correlations, just like the body color of objects in the lion-cat example. The three parts together make up the observed
105 data. Thus even if we exclude all the domain-induced spurious
106 correlations, i.e., entirely remove the effect from \bar{x} , we may still
107 obtain object-induced spurious correlations resulting from \tilde{x} .

108 As Observation 2 shows that excluding only domain-induced spurious correlations is insufficient, we select diverse data batches
109 among the selected domains to help mitigate object-induced spurious correlations. In the level-two-sampling, since we do not have available
110 labels just like domain labels in the level-one-sampling, it is infeasible to utilize DANN again to
111 train a featurizer. So we instead use an ERM model since ERM is prone to taking shortcuts and
112 learning spurious correlations (Zhang et al., 2022). Zhang et al. (2022) also leverages an ERM model
113 to infer the spurious attributes in the unsupervised DG setting. Moreover, since domains attained by
114 the level-one sampling contain diverse data with respect to the domain side, ERM can avert learning
115 domain-induced spurious correlations. Combining these two, the ERM model is prone to relying on
116 object-induced spurious correlations and thus can extract their informative representations. Then
117 a similarity matrix between data batches is constructed with respect to this information. Based on
118 which DPP sampling selects the data batches helping exclude object-induced spurious correlations.

124 2.2 DOMIDiversity Boosted Two-level Sampling

125 Figure 3 shows the sampling procedure of DOMI, a diversity boosted two-level sampling framework.
126 We present the details in Algorithm 1.

Figure 3: Illustration of the sampling procedure of DOMI. The solid arrow indicates the actual sampling flow, while the dotted arrow is used to demonstrate the difference between random sampling and DOMI.

Algorithm 1: Sampling Procedure of DOM

```

Input: The whole training dataset:  $T = \{x_i^d, y_i^d\}_{i=1}^{n_d}$  for  $d \in D$ 
the proportion of domains and batches to be sampled
1 Level-one-sampling
2 Train an invDANN-featurizer  $f$  on a randomly sampled subset  $\bar{d}$ ;
3 for  $d$  in  $D$  do
4    $feat_d \leftarrow \emptyset$ ;
5   for  $i$  from 1 to  $n_d$  do
6      $feat_d \leftarrow feat_d \cup \{x_i^d\}$ ;
7      $feat_d \leftarrow feat_d \cup \{y_i^d\}$ ;
8 Initialize similarity matrix  $L_d = \mathbb{0}_{|D| \times |D|}$ ;
9 for  $d_i$  in  $D$  do
10   for  $d_j$  in  $D$  do
11      $L_d[i, j] \leftarrow \frac{1}{2} (feat_{d_i} - feat_{d_j})^2$ ;
12 Obtain DPP  $L_d$ ;  $\mathbb{P} \sim \text{DPP}(L_d)$ ;  $r = x_i^d, y_i^d, x_{i-1}^d$  for  $d \in D$ ;  $D = L \setminus D$ ;  $\mathbb{P} \leftarrow \mathbb{P} \setminus \mathbb{P}$ ;
13 Level-two-sampling
14 Divide  $r$  into  $R = \{r^b\}_{b=1}^B$  for  $b \in B$ ;
15 Train an ERM featurizer  $f_r$  on  $R$ ;
16 for  $b$  in  $B$  do
17   Compute  $feat_b$  in the same way as computing  $feat_b$  in Level-one-sampling
18 Computing similarity matrix  $L_b$ ;
19 Return  $S = \text{DPP}(L_b)$ ;  $\mathbb{P} \leftarrow \mathbb{P} \setminus \mathbb{P}$ ;

```

3 Experiments

We have investigated the performance of DOM with five backbone DG algorithms on two simulated benchmarks (Rotated MNIST and Rotated Fashion MNIST), which show that DOM can help substantially achieve higher test accuracy. Due to space constraints, and experimental settings, more results and analyses are deferred to Appendix C.1.

Baselines. For each one of the backbone algorithms, we set the baseline as training on domains selected by the random sampling scheme and denoted as $level_0$, compared to the level-one-sampling of DOM and the full version of DOM represented as $level_1$ and $level_2$, respectively. The proportion of minibatches selected in level-two-sampling is a hyperparameter valued from 0 to 1. When equals 1, $level_2$ shrinks to $level_1$.

Results and analysis. Table 1 shows the empirical results and we make the following observations: Strong performance across datasets and algorithms. Considering results on 2 datasets and 5 backbone DG algorithms, $level_1$ gives constant and apparent improvement compared to $level_0$. While $level_2$ may lead to slower growth in accuracy at the initial part of training as shown in Figure 6 because of using a smaller number of minibatches, it keeps outperforming $level_1$ and $level_0$ at later epochs. The

Table 1: Average test accuracy. We repeat the experiment for 5 times on FISH and 20 times on the other algorithms with random seeds.

Dataset	Sampling scheme	DANN	MatchDG	FISH	MMD	CORAL
Rotated MNIST	$level_0$	74.5	81.5	65.2	84.2	85.6
	$level_1$	76.5 ± 2.0	83.6 ± 2.1	66.5 ± 1.3	87.2 ± 3.0	89.2 ± 3.6
	$level_2$	78.6 ± 4.1	84.2 ± 2.7	66.6 ± 1.4	87.7 ± 3.5	89.6 ± 4.0
Rotated Fashion MNIST	$level_0$	40.3	38.2	33.2	39.0	38.7
	$level_1$	42.8 ± 2.5	39.7 ± 1.5	34.5 ± 1.3	41.8 ± 2.8	40.8 ± 2.1
	$level_2$	43.5 ± 3.2	40.7 ± 2.5	35.8 ± 2.6	42.8 ± 3.8	42.1 ± 3.4

gap between test accuracy and maximal accuracy. During training we observe that the test accuracy first rises to the peak value and then begins to decline along with the increase of validation accuracy. This reduction indicates a certain degree of overfitting to spurious correlations. Thus we further record the peak value of the test accuracy in each experiment and denote it as maximal accuracy. The distribution of test accuracy and maximal accuracy on MatchDG under different sampling schemes is shown in Figure 5. While the test accuracy of $level_0$ scatters, that of $level_2$ centers, and $level_2$ shrinks the gap between test accuracy and maximal accuracy.

149 **Ethics statement**

150 This study does not involve any of the following: human subjects, practices to dataset releases,
151 potentially harmful insights, methodologies and applications, potential conflicts of interest and
152 sponsorship, discrimination/bias/fairness concerns, privacy and security issues, legal compliance, and
153 research integrity issues.

154 **Reproducibility statement**

155 To ensure the reproducibility of our empirical results, we present the detailed experimental settings
156 in Appendix C.1 in addition to the main text. Besides, we will further provide the source codes for
157 reproducing results in our paper.

158 **References**

- 159 Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization.
160 arXiv preprint arXiv:1907.02893, 2019.
- 161 Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, Francois
162 Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks.
163 The journal of machine learning research, 17(1):2096–2030, 2016.
- 164 Judy Hoffman, Eric Tzeng, Taesung Park, Jun-Yan Zhu, Phillip Isola, Kate Saenko, Alexei Efros,
165 and Trevor Darrell. Cycada: Cycle-consistent adversarial domain adaptation. *International
166 conference on machine learning*, pp. 1989–1998. Pmlr, 2018.
- 167 Yuanfeng Ji, Lu Zhang, Jiaxiang Wu, Bingzhe Wu, Long-Kai Huang, Tingyang Xu, Yu Rong, Lanqing
168 Li, Jie Ren, Ding Xue, et al. Drugood: Out-of-distribution (ood) dataset curator and benchmark
169 for ai-aided drug discovery—a focus on affinity prediction problems with noise annotations. *arXiv
170 preprint arXiv:2201.09637*, 2022.
- 171 Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Bal-
172 subramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, et al. Wilds: A
173 benchmark of in-the-wild distribution shifts. *International Conference on Machine Learning*
174 pp. 5637–5664. PMLR, 2021.
- 175 Alex Kulesza, Ben Taskar, et al. Determinantal point processes for machine learning. *Foundations
176 and Trends® in Machine Learning*, 5(2–3):123–286, 2012.
- 177 Haoliang Li, Sinno Jialin Pan, Shiqi Wang, and Alex C Kot. Domain generalization with adversarial
178 feature learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*
179 pp. 5400–5409, 2018.
- 180 Jiashuo Liu, Zheyuan Hu, Peng Cui, Bo Li, and Zheyuan Shen. Heterogeneous risk minimization. In
181 *International Conference on Machine Learning*, pp. 6804–6814. PMLR, 2021.
- 182 Mingsheng Long, Yue Cao, Jianmin Wang, and Michael Jordan. Learning transferable features with
183 deep adaptation networks. *International conference on machine learning*, pp. 97–105. PMLR,
184 2015.
- 185 Divyat Mahajan, Shruti Tople, and Amit Sharma. Domain generalization using causal matching. In
186 *International Conference on Machine Learning*, pp. 7313–7324. PMLR, 2021.
- 187 Alfred Müller. Integral probability metrics and their generating classes of functions. *Advances in
188 Applied Probability* 29(2):429–443, 1997.
- 189 Vihari Piratla, Praneeth Netrapalli, and Sunita Sarawagi. Efficient domain generalization via common-
190 specific low-rank decomposition. In *International Conference on Machine Learning*, pp. 7728–
191 7738. PMLR, 2020.

- 192 Yuge Shi, Jeffrey Seely, Philip HS Torr, N Siddharth, Awni Hannun, Nicolas Usunier, and Gabriel
193 Synnaeve. Gradient matching for domain generalization. arXiv preprint arXiv:2104.09937, 2021.
- 194 Baochen Sun and Kate Saenko. Deep coral: Correlation alignment for deep domain adaptation. In
195 European conference on computer vision. pp. 443–450. Springer, 2016.
- 196 Michael Zhang, Nimit S Sohoni, Hongyang R Zhang, Chelsea Finn, and Christopher Ré. Correct-n-
197 contrast: A contrastive approach for improving robustness to spurious correlations. arXiv preprint
198 arXiv:2203.01517, 2022.
- 199 Han Zhao, Shanghang Zhang, Guanhang Wu, José MF Moura, Joao P Costeira, and Geoffrey J Gordon.
200 Adversarial multiple source domain adaptation. Advances in neural information processing systems
201 31, 2018.
- 202 Han Zhao, Remi Tachet Des Combes, Kun Zhang, and Geoffrey Gordon. On learning invariant
203 representations for domain adaptation. International Conference on Machine Learning, pp.
204 7523–7532. PMLR, 2019.

Appendix of DOMI

205

206 Contents

207 **A Theoretical analysis** 7

208 **B The Simulated Dataset** 8

209 **C Experimental Details** 8

210 **C.1 Settings and results** 8

211 **C.2 Experiments on iwildcam** 11

212 **D How can spurious correlations occur in the two datasets?** 11

213 A Theoretical analysis

214 Preliminaries. Consider the universal set of domains \mathcal{D} , where each domain $d \in \mathcal{D}$ corresponds
 215 to a distribution P_d over $X \times Y$, with X being the space of inputs and that of outputs. Our goal
 216 is to find a predictor $f: X \rightarrow Y$ while we can only access the domains \mathcal{D}_Π where $\mathcal{D}_\Pi \subseteq \mathcal{D}$.
 217 We measure the quality of a prediction with a loss function $\ell: Y \rightarrow \mathbb{R}_0^+$; and the quality of a
 218 predictor by its population loss on domain $d \in \mathcal{D}$, given by $L_d(f) = \mathbb{E}_{x,y \sim P_d} \ell(f(x); y)$.

219 Definition 1 (Correlation). We define a correlation as a predictor $f: X \rightarrow Z$ where $X \rightarrow Z$ is a
 220 data representation and $Z \rightarrow Y$ is a classifier. The causal correlation satisfies: elicits
 221 an invariant predictor (Arjovsky et al., 2019) and: f is simultaneously optimal for all domains, i.e.,
 222 $f = \arg\min_{f: X \rightarrow Y} L_d(f) \quad \forall d \in \mathcal{D}$.

223 Notably, Definition 1 requires that f and g are unrestricted in the space of all (measurable) functions.
 224 However, we learn f and g being restricted to only access domains \mathcal{D}_Π , a small subset of \mathcal{D} . For
 225 this to be feasible, it is natural to add a restriction that f and g belong to suitable classes of
 226 functions mapping $X \rightarrow Z$ and W of functions mapping $Z \rightarrow Y$.

227 Assumption 1. $\arg\min_{f: X \rightarrow Y} L_d(f) = \arg\min_{f \in \mathcal{F}} L_d(f)$ & x ; where \mathcal{F} is a constant.

228 Definition 2. Consider a domain set \mathcal{D}_s , on which the set of invariant predictors, \mathcal{D}_s , is the set
 229 of all predictors f satisfies following:

- 230 a $f \in \mathcal{F}$ with $f = \arg\min_{f \in \mathcal{F}} L_d(f)$;
- 231 a for all $d \in \mathcal{D}_s, f = \arg\min_{f \in \mathcal{F}} L_d(f)$.

232 Lemma A.1. Based on Definition 1 and Definition 2, we can trivially derive: for any nonempty set
 233 $\mathcal{D} \subseteq \mathcal{D}; f = \arg\min_{f \in \mathcal{F}} L_d(f)$.

234 Definition 3 (Diversity). We use Integral Probability Metric (Müller, 1997) to measure the diversity
 235 between domains. For domain $d \in \mathcal{D}$, the diversity is defined as:

$$\text{Div}(P_d; P_d) = \text{Div}(P_d; P_d; \mathcal{G}) = \sup_{g \in \mathcal{G}} |\mathbb{E}_{P_d} g(x; y) - \mathbb{E}_{P_d} g(x; y)|$$

236 Where \mathcal{G} is a class of bounded functions. When we let $y = f(x)$ and $\mathcal{G} = \mathcal{F}$, the
 237 diversity is:

$$\text{Div}(P_d; P_d) = \text{Div}(P_d; P_d; \mathcal{F}) = \sup_{f \in \mathcal{F}} |\mathbb{E}_{P_d} \ell(f(x); y) - \mathbb{E}_{P_d} \ell(f(x); y)|$$

$$\text{Div}(P_d; P_d; \mathcal{F}) = \sup_{f \in \mathcal{F}} |L_d(f) - L_d(f)|$$

238 Consider we have a domain $\mathcal{D}_k = \{d_1; d_2; \dots; d_k\}$ and the corresponding $D_k = \{f_1; f_2; \dots; f_m\}$.
 239 And now we get one more domain d_{k+1} to form D_{k+1} . According to Lemma A.1, the causal
 240 correlation $f = \arg\min_{f \in \mathcal{F}} L_d(f)$, so an informative domain d_{k+1} which helps exclude spurious correlations
 241 leads to $f = \arg\min_{f \in \mathcal{F}} L_d(f)$.

242 Proposition 1 (Diverse domains help exclude spurious correlations). d_{k-1} satisfies that:
 243 $\max_{d \in D_k; f_i \in I_{D_k}} \text{Div } d_{k-1}; d \perp L_d f_i$ & , then $I_{D_k} \perp I_{D_{k-1}}$.

244 Proof. Without loss of generality, we first conduct analysis on I_{D_k} . For f_t :

$$\max_{d \in D_k} \mathbb{1}_{L_{d_{k-1}} f_t \perp L_d f_t} \perp L_d f_t \text{ \& } \max_{d \in D_k} \text{Div } d_{k-1}; d \perp L_d f_t$$

$$\max_{d \in D_k} \text{Div } d_{k-1}; d \perp L_d f_t \text{ \& } \max_{d \in D_k; f_i \in I_{D_k}} \text{Div } d_{k-1}; d \perp L_d f_i \text{ \& }$$

When

$$a \perp L_{d_{k-1}} f_t \perp L_d f_t \text{ \& } 0$$

$$L_{d_{k-1}} f_t \perp L_d f_t \text{ \& }$$

$$a \perp L_{d_{k-1}} f_t \perp L_d f_t \text{ \& } 0$$

$$L_{d_{k-1}} f_t \perp \max_{d \in D_k} L_{d_{k-1}} f_t \perp L_d f_t \perp L_d f_t \perp \max_{d \in D_k} \mathbb{1}_{L_{d_{k-1}} f_t \perp L_d f_t} \perp L_d f_t \text{ \& }$$

245 $L_{d_{k-1}} f_t \text{ \& } ,$ we get $f_t \perp I_{D_{k-1}}$ for any $t \in \{1, 2, \dots, m\}$, thus $I_{D_k} \perp I_{D_{k-1}}$ □

246 B The Simulated Dataset

Table 2: The simulated dataset of the toy example. From these 12 data points, we sample 6 for training.

	D ₁	D ₂	D ₃	D ₄	D ₅	D ₆	D ₇	D ₈	D ₉	D ₁₀	D ₁₁	D ₁₂
X ₁	0	0	0	0	0	0	0	0	0	1	1	1
X ₂	0	0	0	0	0	0	1	1	1	0	1	1
X ₃	0	0	0	0	1	1	1	1	1	1	1	1
X ₄	0	0	0	1	0	1	1	1	0	0	1	1
Y	0	0	0	0	0	0	1	1	1	1	1	1

247 C Experimental Details

248 C.1 Settings and results

249 **Datasets.** To satisfy the setting of a large number of domains, we extend the original simulated
 250 benchmarks on MNIST and Fashion MNIST by Piratla et al. (2020) from rotating images by
 251 75° in intervals of 15° to intervals of 1° in the training set, i.e., 61 domains in total. And we get test
 252 accuracy on the test set which rotates images either 90°. Moreover, while the original datasets
 253 rotate the same images for different degrees, we extend them to fit the real cases in DG tasks. We
 254 generate indices using different random seeds to select images from MNIST and Fashion MNIST
 255 for each domain before rotating. Appendix D gives examples to show how spurious correlations can
 256 occur in the two datasets.

257 **Backbones.** We take MatchDG (Mahajan et al., 2021), FISH (Shi et al., 2021), CORAL (Sun
 258 & Saenko, 2016), MMD (Li et al., 2018) and DANN (Ganin et al., 2016) as backbone algorithms.
 259 The former four algorithms work by pairwise domain operations, leading to O^2 computational
 260 complexity with n domains and thus prohibitive to be scaled to DG tasks with multitudinous domains.
 261 It is essential for them to sample the most informative domains. We further incorporate DANN as
 262 one of the backbone algorithms since DANN can not only efficiently select domains by its first level
 263 of sampling but can help deal with circumstances where each domain contains massive data by the
 264 second level of sampling.

Hyperparameters. For DANN, the training epochs are set to be 50. MatchDG is a two-phase method, and in our experiment, we set 30 epochs of training for phase 1 and 25 epochs for phase 2. While $level_1$ gets higher accuracy on Rotated MNIST and $level_2$ shows better performance on Fashion MNIST, they all outperform $level_0$, i.e., randomly sampling. The training epochs of FISH are set to be 5. Each epoch contains 300 iterations and we observe test accuracy every 30 iterations. And in Figure 6 we slightly abuse epoch to mean the time we obtain test accuracy. Unlike MatchDG and DANN, sh needs to sample domains in each iteration instead of training on one list of domains. Sampling domains in each iteration will result in great computational overhead compared to randomly sampling. Thus we just sample 30 domain lists containing diverse domains using level-one-sampling of DOM and repeatedly train the model on these domain lists (one list for one iteration). As for $level_2$, we further utilize level-two sampling to sample data batches of each domain in the domain lists for training. The former 3 DG algorithms utilize SGD optimizer with a learning rate of 0.01, weight decay 5×10^{-4} , and momentum 0.9. The training epochs of MMD and CORAL are set as 30. These two algorithms leverage Adam optimizer with a learning rate of 0.001 and weight decay of 0. All ve algorithms use the Resnet18 model. Within each backbone algorithm, we keep factors including learning rate, batch size, choice of the optimizer, and model architecture the same for $level_0$, $level_1$ and $level_2$ to highlight the effect of different sampling schemes. It's worth noting that we do no comparison between the backbone algorithms since we do not conduct meticulous hyperparameter tuning for them.

Model selection. During training, we use a validation set to measure the model's performance. The test accuracy of the model is updated after an epoch if it shows better validating performance. That is, we save the model with the highest validation accuracy after the training procedure, obtain its test accuracy, and report results. For Rotated MNIST and Rotated Fashion MNIST, data from only source domains (rotation degree is from 15° to 75°) are used to form the validation set.

Empirical results Figure 5 show test accuracy and maximal accuracy among 20 times of repeated experiments with random seeds leveraging different sampling levels on Rotated Fashion MNIST and Rotated MNIST. Among training epochs, the test accuracy rises to the peak value and then declines along with the increase of validation accuracy. In this gure, maximal accuracy represents the peak value. Each tiny circle represents one time of the experiment, of which the vertical location corresponds to the accuracy value. The horizontal line inside each box indicates the mean value.

The choice of α . A smaller α helps efficiently mitigate strong object-induced spurious correlations and speed up training, but when the impact from the object side is weak, a small α leads to a waste of training data. In the experiment, we observe that a relatively small α is more beneficial for Rotated Fashion MNIST while a large α works better on Rotated MNIST. Figure 4 shows the results of different α .

(a) Rotated Fashion MNIST (b) Rotated MNIST
 Figure 4: Average test accuracy of 20 experiments with random seeds during 50 epochs with different α on Rotated Fashion MNIST and Rotated MNIST of DANN. 1:0 corresponds to DOM with only level one.

(a) MatchDG on Rotated Fashion MNIST

(b) MatchDG on Rotated MNIST

(c) DANN on Rotated MNIST

(d) DANN on Rotated Fashion MNIST

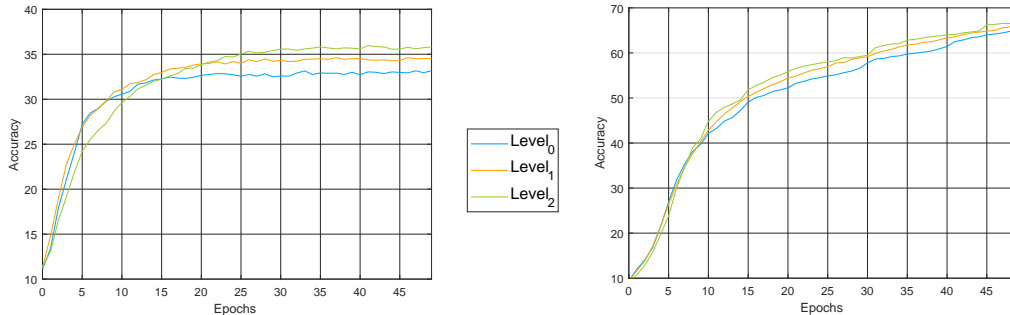
(e) CORAL on Rotated MNIST

(f) CORAL on Rotated Fashion MNIST

(g) MMD on Rotated MNIST

(h) MMD on Rotated Fashion MNIST

Figure 5: Boxplot of test accuracy and maximal accuracy among 20 repeated experiments with random seeds leveraging different sampling levels on Rotated Fashion MNIST and Rotated MNIST. Among training epochs, the test accuracy rises to the peak value and then declines with the increase of validation accuracy. In this figure, maximal accuracy represents the peak value. Each tiny circle represents one time of the experiment, of which the vertical location corresponds to the accuracy value. The horizontal line inside each box indicates the mean value.



(a) Rotated Fashion MNIST

(b) Rotated MNIST

Figure 6: Average test accuracy of 5 experiments with random seeds during 50 epochs under different sampling schemes of FISH.

300 C.2 Experiments on iwildcam

301 WILDS (Koh et al., 2021) is a curated collection of benchmark datasets representing distribution
 302 shifts faced in the wild. As one dataset in WILDS, iwildcam contains photos of wild animals and 324
 303 different camera traps are taken as domains. The data of iwildcam is extremely unbalanced, while part
 304 of the domains contains less than 20 photos, some domains contain over 2000 ones. In the original
 305 experiments of Shi et al. (2021), iwildcam is divided into batches in each domain. FISH samples 10
 306 batches from different domains for training in each iteration. The sampling probability of one batch
 307 in a domain is proportional to the number of batches left in this domain. This sampling scheme is
 308 taken as $level_0$ here and we refer to the result of (Shi et al., 2021). In each iteration, $level_1$ samples
 309 10 batches based on DPP using invDANN, $level_2$ first samples 10 batches in the level-one-sampling
 310 and among them selects 6 batches in the level-two-sampling. Under the same setting in the original
 311 experiments, the results on iwildcam of FISH are shown in Table 3 .

Table 3: Macro F1 score of FISH on iwildcam under three sampling schemes

	$level_0$	$level_1$	$level_2$
Iwildcam	22.0	22.8	23.4

312

313 Although DOMI gets a higher Macro F1 score, it leads to a much larger computational overhead since
 314 it needs to do sampling in each iteration. Moreover, for DANN and MatchDG, Macro F1 of diverse
 315 domains may be significantly lower than randomly sampled domains because of the unbalanced
 316 data, i.e., the diverse domains may contain much fewer data compared to the randomly sampled
 317 domains. It would be significant future work to tackle the issues of extremely imbalanced data and
 318 computational overhead for algorithms that need to do sampling for multi-times.

319 D How can spurious correlations occur in the two datasets?

320 It’s much easier to differentiate the rotation degree than to discriminate the objects. This can be
 321 empirically verified since it only needs about 30 epochs for a model to achieve over 98% validation
 322 accuracy of classifying 61 different degrees while 50 epochs to achieve no more than 97% and
 323 88% validation accuracy of classifying 10 different objects on rotated MNIST and Fashion MNIST,
 324 respectively. Thus if a certain class label is closely associated with a certain rotation degree in the
 325 training set, recognizing objects by actually recognizing the rotation degree can be a shortcut and
 326 domain-induced spurious correlation, just like classifying cats and lions using the background in the
 327 toy example. As for object-induced spurious correlation, on rotated MNIST, the handwriting is the
 328 feature of the object, however, it can also be the spurious correlation. For example, in Figure 7, let’s
 329 focus on the number “1” and “7”. After training on Figure 7a, can the model correctly recognize “1”
 330 in Figure 7b instead of wrongly taking it as “7”?

