

SPAWRIOUS: A BENCHMARK FOR FINE CONTROL OF SPURIOUS CORRELATION BIASES

Anonymous authors

Paper under double-blind review

ABSTRACT

The problem of spurious correlations (SCs) arises when a classifier relies on non-predictive features that happen to be correlated with the labels in the training data. Previous SC benchmark datasets suffer from varying issues, e.g., over-saturation or only containing one-to-one (O2O) SCs, but no many-to-many (M2M) SCs arising between groups of spurious attributes and classes. In this paper, we present Spawrious- $\{O2O, M2M\}$ - $\{Easy, Medium, Hard\}$, an image classification benchmark suite containing spurious correlations between classes and backgrounds. We employ a text-to-image model to generate photo-realistic images and an image captioning model to filter out unsuitable ones. The resulting dataset is of high quality and contains approximately 152k images. Our experimental results demonstrate that state-of-the-art group robustness methods struggle with Spawrious.

1 INTRODUCTION

To make models more robust to unseen test distributions, mitigating a classifier’s reliance on spurious, non-causal features that are not essential to the true label has attracted lots of research interest (Sagawa et al., 2019a; Arjovsky et al., 2019; Kaddour et al., 2022b; Izmailov et al., 2022). For example, classifiers trained on ImageNet (Deng et al., 2009) have been shown to rely on backgrounds (Xiao et al., 2020; Singla & Feizi, 2022; Neuhaus et al., 2022), which are spuriously correlated with class labels but, by definition, not predictive of them.

Recent work has focused substantially on developing new methods for addressing the spurious correlations (SCs) problem (Kaddour et al., 2022b), yet, studying and addressing the limitations of existing benchmarks remains underexplored. For example, the *Waterbirds* (Sagawa et al., 2019a), and *CelebA hair color* (Liu et al., 2015) benchmarks remain among the most used benchmarks for the SC problem; yet, GroupDRO (Sagawa et al., 2019a) achieves 90.5% worst-group accuracy using group adjusted data with a ResNet50 pretrained on ImageNet.

Another limitation of existing benchmarks is their sole focus on overly simplistic one-to-one (O2O) spurious correlations, where one spurious attribute correlates with one label. However, in reality, we often face *many-to-many* (M2M) spurious correlations across groups of classes and backgrounds, which we formally introduce in this work.

While some benchmarks include multiple training environments with varying correlations (Koh et al., 2021), they do not test classification performance on reversed correlations during test time. Such M2M-SCs are *not* an aggregation of O2O-SCs and cannot be expressed or decomposed in the form of the latter; they contain qualitatively different spurious structures, as shown in Figure 2. To our knowledge, this work is the first to conceptualize and instantiate M2M-SCs in image classification problems.

Contributions We introduce *Spawrious*- $\{O2O, M2M\}$ - $\{Easy, Medium, Hard\}$, a suite of image classification datasets with O2O and M2M spurious correlations and three difficulty levels each. Recent work (Wiles et al., 2022; Lynch et al., 2022; Vendrow et al., 2023) has demonstrated a proof-of-concept to effectively discover spurious correlation failure cases in classifiers by leveraging off-the-shelf, large-scale, image-to-text models trained on vast amounts of data. Here, we take this view to the extreme and generate a novel benchmark with 152,064 images of resolution 224×224 , specifically targeted at the probing of classifiers’ reliance on spurious correlations.

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

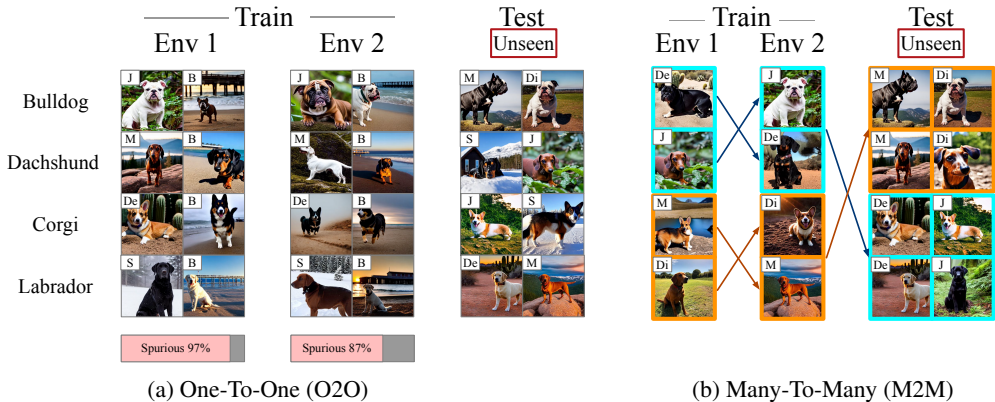


Figure 1: **Spawrious Challenges:** Letters on the images denote the background, and the bottom bar in Figure 1a indicates each class’s proportion of the spurious background. In the O2O challenge, each class associates with a background during training, while the test data contains unseen combinations of class-background pairs. In the M2M challenge, a group of classes correlates with a group of backgrounds during training, but this correlation is reversed in the test data.

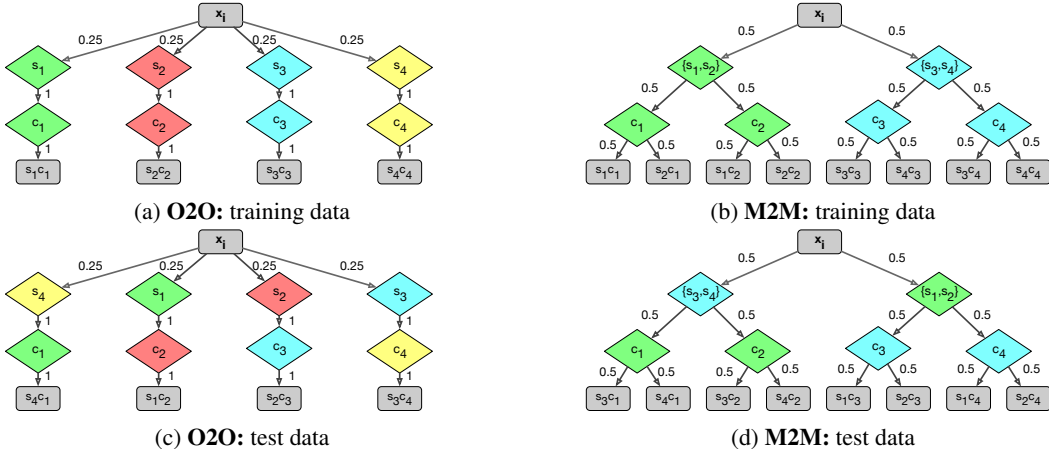


Figure 2: **Data distributions for our challenges:** x_i is a random image sampled, each s_i is a spurious attribute, and each c_i is a class label. The edges indicate the probability that the sample x_i has a given property, conditional on previous steps in the tree. The leaf nodes indicate the possible attribute-class combinations in the distribution. The colors emphasize the distribution shift in the test data.

Our experimental results demonstrate that state-of-the-art methods struggle with Spawrious, most notably on the *Hard*-splits with $< 73\%$ accuracy using ResNet50 pretrained on ImageNet.

2 BENCHMARK DESIDERATA

2.1 SIX DESIDERATA

1. Photo-realism, unlike datasets containing cartoon/sketch images (Gulrajani & Lopez-Paz, 2021) or image corruptions (Hendrycks & Dietterich, 2019), which are known to conflict with current backbone network architectures (Geirhos et al., 2018a;b; Hermann et al., 2020), possibly confounding the evaluation of OOD algorithms. **2. Non-binary classification problem**, to minimize accidentally correct classifications achieved by chance. **3. Inter-class homogeneity and intra-class heterogeneity**, i.e., low variability *between* and high variability *within* classes, to minimize the margins of the decision boundaries inside the data manifold (Murphy, 2022). This desideratum ensures that the classification problem is non-trivial. **4. High-fidelity backgrounds** including complex features

to reflect realistic conditions typically faced in the wild instead of monotone or entirely removed backgrounds (Xiao et al., 2020). **5. Access to multiple training environments**, i.e., the conditions of the *Domain Generalization* problem (Gulrajani & Lopez-Paz, 2021), which allow us to learn domain invariances, such that classifiers can perform well in novel test domains. **6. Multiple difficulty levels**, so future work can study cost trade-offs. For example, one may budget higher computational costs for methods succeeding on difficult datasets than those that succeed only on easy ones.

2.2 SPURIOUS CORRELATIONS (ONE-TO-ONE)

Here, we provide some intuition and discuss the conditions for a (one-to-one) spurious correlation (SC). We define a correlated, non-causal feature as a feature that frequently occurs with a class but does not cause the appearance of the class (nor vice versa). We abuse the term “correlated” as it is commonly used by previous work, but we consider non-linear relationships between two random variables too. Further, we call correlated features *spurious* if the classifier perceives them as a feature of the correlated class.

Next, we want to define a *challenge* that allows us to evaluate a classifier’s harmful reliance on spurious features. Spurious features are not always harmful; even humans use context information to make decisions (Geirhos et al., 2020). However, a spurious feature becomes harmful if it alone is sufficient to trigger the prediction of a particular class without the class object being present in the image (Neuhauss et al., 2022).

To evaluate a classifier w.r.t. such harmful predictions, we evaluate its performance when the spurious correlations are reverted. The simplest setting is when a positive/negative correlation exists between one background variable and one label in the training/test environment.

O2O-SC Challenge

Let $p(\mathbf{X}, S, C)$ be a distribution over images $\mathbf{X} \in \mathbb{R}^D$, spurious attributes $S \in \mathcal{S} = \{s_1, \dots, s_K\}$, and labels $C \in \mathcal{C} = \{c_1, \dots, c_P\}$. Given $\hat{p}_{\text{data}} \neq p_{\text{test}}$, and $K = P$ it holds that for $i \in [K]$,

$$\text{corr}_{\hat{p}_{\text{data}}}(\mathbb{1}(S = s_i), \mathbb{1}(C = c_i)) > 0, \text{corr}_{p_{\text{test}}}(\mathbb{1}(S = s_i), \mathbb{1}(C = c_i)) < 0. \quad (1)$$

where the indicator function $\mathbb{1}(X = x)$ is non-zero when the *variable* X equals the *value* x .

Figure 1a illustrates the one-to-one (O2O) SC, in which pair-wise SCs between spurious features S and labels C exist within training environments, which then differ in the test environment.

2.3 MANY-TO-MANY SPURIOUS CORRELATIONS

Figure 2 shows an example of how to construct M2M-SCs, which contain richer spurious structures, following an *hierarchy* of the class groups correlating with spurious attribute groups. As we will see later in Section 3.2, the data-generating processes we instantiate for each challenge differ qualitatively.

M2M-SC Challenge

Consider $p(\mathbf{X}, S, C)$ defined in the O2O-SC Challenge. We further assume the existence of partitions $\mathcal{S} = \mathcal{S}_1 \dot{\cup} \mathcal{S}_2$ and $\mathcal{C} = \mathcal{C}_1 \dot{\cup} \mathcal{C}_2$. Given $\hat{p}_{\text{data}}, p_{\text{test}}$, it holds that for $j \in \{1, 2\}$

$$\text{corr}_{\hat{p}_{\text{data}}}(\mathbb{1}(S \in \mathcal{S}_j), \mathbb{1}(C \in \mathcal{C}_j)) = 1, \text{corr}_{p_{\text{test}}}(\mathbb{1}(S \in \mathcal{S}_j), \mathbb{1}(C \in \mathcal{C}_j)) = -1. \quad (2)$$

3 THE SPAWRIOUS CHALLENGE

3.1 DATASET CONSTRUCTION

We instantiate the desiderata introduced in Section 2 by presenting *Spawrious*, a synthetic image classification dataset containing images of four dog breeds (classes) in six background locations (spurious attributes). Figure 3 summarizes the dataset construction pipeline, which we now discuss

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

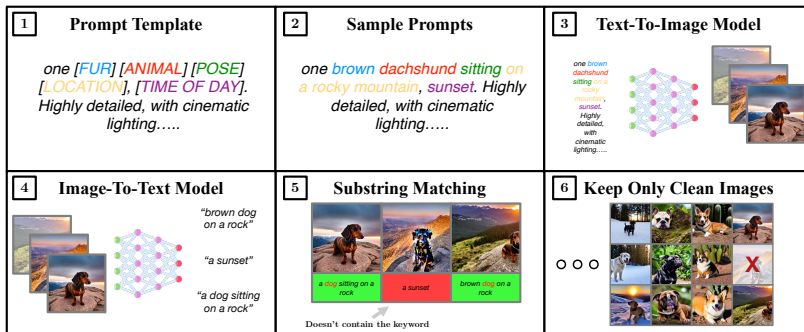


Figure 3: **Spawrious Pipeline**: We leverage text-to-image models for generation (Steps 1-3) and image-to-text models for cleaning of bad images (Steps 4-6). Details in Section 3.1 and Appendix H.

in more detail. The main idea is to leverage recently proposed text-to-image models (Rombach et al., 2022) for photo-realistic image generation and image-to-text models (NLP Connect, 2022) for filtering out low-quality images. We address potential ethical concerns that may arise from using a generative model to construct this dataset in Appendix D.

A **prompt template** allows us to define high-level factors of variation. We then **sample prompts** by filling in randomly sampled values for these high-level factors. The **text-to-image model** generates images given a sampled prompt; we use *Stable Diffusion v1.4* (Rombach et al., 2022). We pass the raw, generated images to an **image-to-text (I2T) model** to extract a concise description; here, we use the ViT-GPT2 image captioning model (NLP Connect, 2022). We perform a form of **substring matching** by checking whether important keywords are present in the caption, e.g., “dog”. This step avoids including images without class objects, which we sometimes observed due to the T2T model ignoring parts of the input prompt. We **keep only “clean” images** whose captions include important keywords. More details on this pipeline and possible failures are discussed in Appendix H, as well as a measure of the accuracy of the prompt-image alignment in Appendix I.

3.2 SATISFYING BENCHMARK DESIDERATA

To ensure **photorealism**, we generate images using *Stable Diffusion v1.4* (Rombach et al., 2022), trained on a large-scale real-world image dataset (Schuhmann et al., 2022), while carefully filtering out images without detectable class objects. We construct a 4-way classification problem to reduce the probability of accidentally correct classifications compared to a **binary classification problem** (e.g., CelebA hair color prediction or Waterbirds). Next, we chose dog breeds to reduce **inter-class variance**, inspired by the difference in classification difficulty between Imagenette (easily classified objects) (Howard, 2019a), and ImageWoof (Howard, 2019b) (dog breeds), two datasets based on subsets of ImageNet (Deng et al., 2009). We increase **intra-class variance** by adding animal poses to the prompt template.

We add “[location] [time of day]” variables to the prompt template to ensure **diverse backgrounds**, and select six combinations after careful experimentation with dozens of possible combinations, abandoning over-simplistic ones. Our final prompt template takes the form “one [fur] [animal] [pose] [location], [time of day]. highly detailed, with cinematic lighting, 4k resolution, beautiful composition, hyperrealistic, trending, cinematic, masterpiece, close up”, and there are 72 possible combinations. The variables [location]/[animal] correspond to spurious backgrounds/labels for a specific background-class combination. The other variables take the following values: “fur: black, brown, white, [empty]; pose: sitting, running, [empty]; time of day: pale sunrise, sunset, rainy day, foggy day, bright sunny day, bright sunny day”.

To construct **multiple training environments**, we randomly sample from a set of background-class combinations, which we further group by their **difficulty level into easy, medium, and hard**. We construct two datasets for each SC type with 3, 168 images per background-class combination, thus $2 \text{ SC types} \times 4 \text{ environments} \times 6 \text{ difficulties} \times 3, 168 = 152, 064$ images in total.

Class	Train Env 1	Train Env 2	Test	Train Env 1	Train Env 2	Test	Train Env 1	Train Env 2	Test
	O2O-Easy			O2O-Medium			O2O-Hard		
Bulldog	97% De 3% B	87% De 13% B	100% Di	97% M 3% De	87% M 13% De	100% J	97% J 3% B	87% J 13% B	100% M
Dachshund	97% J 3% B	87% J 13% B	100% S	97% B 3% De	87% B 13% De	100% Di	97% M 3% B	87% M 13% B	100% S
Labrador	97% Di 3% B	87% Di 13% B	100% De	97% Di 3% De	87% Di 13% De	100% B	97% S 3% B	87% S 13% B	100% De
Corgi	97% S 3% B	87% S 13% B	100% J	97% J 3% De	87% J 13% De	100% S	97% De 3% B	87% De 13% B	100% J
	M2M-Easy			M2M-Medium			M2M-Hard		
Bulldog	100% Di	100% J	50% S 50% B	100% De	100% M	50% Di 50% J	100% B	100% S	50% De 50% M
Dachshund	100% J	100% Di	50% S 50% B	100% M	100% De	50% Di 50% J	100% B	100% S	50% De 50% M
Labrador	100% S	100% B	50% Di 50% J	100% Di	100% J	50% De 50% M	100% M	100% De	50% B 50% S
Corgi	100% B	100% S	50% Di 50% J	100% J	100% Di	50% De 50% M	100% M	100% De	50% B 50% S

Table 1: Proportions of Spurious Backgrounds By Class and Environment. Backgrounds include: Beach (B), Desert (De), Dirt (Di), Jungle (J), Mountain (M), Snow (S).

Method	One-To-One SC			Many-To-Many SC			Average
	Easy	Medium	Hard	Easy	Medium	Hard	
ERM (Vapnik, 1991)	77.49%±0.05	76.60%±0.02	71.32%±0.09	83.80%±0.01	53.05%±0.03	58.70%±0.04	70.16%
GroupDRO (Sagawa et al., 2019a)	80.58%±0.74	75.96%±2.18	76.99%±2.60	79.96%±2.79	61.01%±4.64	60.86%±1.71	72.56%
IRM (Arjovsky et al., 2019)	75.45%±2.57	76.39%±2.22	74.90%±1.27	76.15%±2.83	67.82%±4.39	60.93%±1.09	71.94%
CORAL (Sun & Saenko, 2016)	89.66%±1.23	81.05%±1.20	79.65%±1.82	81.26%±1.61	65.18%±4.85	67.97%±0.91	77.46%
CausIRL (Chevalley et al., 2022)	89.32%±1.20	78.64%±0.62	80.40%±1.32	85.76%±1.02	63.15%±2.98	68.93%±0.28	77.20%
MMD-AAE (Li et al., 2018)	78.81%±0.02	75.33%±0.03	72.66%±0.01	80.55%±0.02	59.43%±0.04	54.39%±0.05	70.20%
Fish (Shi et al., 2021)	77.51%±1.58	77.72%±2.82	74.73%±2.40	81.60%±3.44	63.03%±1.96	58.94%±2.56	72.26%
VRex (Krueger et al., 2020)	84.69%±1.69	77.56%±0.62	75.41%±2.67	81.22%±1.25	54.28%±5.42	59.21%±5.08	72.06%
W2D (Huang et al., 2022)	81.94%±1.03	76.74%±0.70	76.84%±1.32	80.80%±2.24	62.82%±2.23	61.89%±2.71	73.50%
JTT (Zheran Liu et al., 2021)	90.24%±3.09	87.28%±0.91	87.41%±0.99	79.23%±1.83	60.56%±5.55	57.58%±3.86	77.05%
Mixup (Xu et al., 2019) // random shuffle	88.48%±0.74	82.75%±3.12	75.75%±1.16	89.61%±0.66	77.23%±0.97	71.21%±2.33	80.84%
Mixup // LISA (Yao et al., 2022)	88.64%±0.51	80.83%±1.33	72.54%±1.07	87.24%±2.51	71.78%±0.31	72.97%±4.23	79.00%

Table 2: Results for Spurious-{O2O,M2M}-{Easy, Medium, Hard} using ImageNet-pretrained ResNet-50: JTT (Zheran Liu et al., 2021) performs the best across the O2O challenges, while Mixup methods (Xu et al., 2019) perform best across M2M challenges and overall attain the highest average.

O2O-SC Challenge We select combinations such that each class is observed with two backgrounds, spurious b^{sp} and generic b^{ge} . For all images with class label c_i in the training data, $\mu\%$ of them have the spurious background b_i^{sp} and $(100 - \mu)\%$ of them have the generic background b^{ge} . Importantly, each spurious background is observed with only one class ($\hat{p}_{data}(b_i^{sp} | c_j) = 1$ if $i = j$ and 0 for $i \neq j$), while the generic background is observed for all classes with equal proportion. We train on two separate environments (with distinct data) that differ in their μ values. Thus, the change in this proportion should serve as a signal to a robustness-motivated optimization algorithm (e.g. IRM (Arjovsky et al., 2019), GroupDRO (Sagawa et al., 2019a) etc.) that the correlation is spurious.

M2M-SC Challenge First, we construct disjoint background and class groups $\mathcal{S}_1, \mathcal{S}_2, \mathcal{C}_1, \mathcal{C}_2$, each with two elements. Then, we select background-class combinations for the training data such that for each class $c \in \mathcal{C}_i$, we pick a combination (s, b) for each $s \in \mathcal{S}_i$. Second, we introduce two environments as shown in Figure 1b.

4 EXPERIMENTS

We fine-tune a ResNet50 (He et al., 2016) model pre-trained on ImageNet, following previous work on domain generalization (Dou et al., 2019; Li et al., 2019; Gulrajani & Lopez-Paz, 2021). Given the size of our dataset, in preliminary experiments, we also tried training a ResNet50 from scratch, which consistently led to worse results. See Appendix E for analysis on the effect of ImageNet pretraining. We use various popular OOD methods, as listed below.

We find that JTT performs the best on the O2O challenges while being one of the worst methods on the M2M challenges. Within the M2M challenge, we find Mixup to perform the best, for both random shuffle and LISA, and overall Mixup attains the best average. This result contributes to the debate whether, for a fixed architecture, most robustness methods perform about the same (Gulrajani & Lopez-Paz, 2021) or not (Wiles et al., 2021). The performances of most methods get consistently worse as the challenge becomes harder. Most often, the data splits of our newly formalized M2M-SC are significantly more challenging than the O2O splits, most notably $M2M\{-Hard, Medium\}$. We conjecture that there is a strong need for new methods targeting such. {ERM, GroupDRO} and {CORAL, CausIRL} perform about the same, despite much different robustness regularization. All methods consistently achieve 98-99% in-distribution test performance (not shown in Table 3 to save

space) despite differences in OOD performance. ERM performs worst on average for the ResNet50 set of results.

5 RELATED WORK

We summarized related benchmarks in Appendix A. Further, we outline some works closest to ours here.

Out-of-distribution Generalization approaches involve training a model simultaneously on multiple related but different domains, exploiting additional environment index labels in the training data (Ben-David et al., 2010; Blanchard et al., 2011; Muandet et al., 2013; Arjovsky et al., 2019), which our benchmark provides too. In order to design effective training losses, approaches may optimize the loss on the worst performing environment (Sagawa et al., 2019a), or enforce an invariance constraint, such as on the features (Sun & Saenko, 2016; Arjovsky et al., 2019; Chevalley et al., 2022) or on the gradients (Rame et al., 2022a). We discuss the methods we applied to our benchmark in Appendix C.

Spurious Correlations have a long history in mathematical statistics (Pearson, 1897; Simon, 1954) and recently entered the machine learning discourse Sagawa et al. (2019b; 2020); Izmailov et al. (2022). They have been detected in common image classification settings via the usage of saliency maps (Moayeri et al., 2022a; Singla & Feizi, 2022). We use saliency maps to validate that an ERM model trained on Spawrious learned dependence on the spurious background feature in Appendix G.

6 LIMITATIONS AND FUTURE WORK

The main limitations of our work have to do with how flexible the dataset can be. Spurious correlations can include **non-background** spurious attributes which currently are not covered. For example, Neuhaus et al. (2022) find that in the ImageNet (Deng et al., 2009) dataset, the class “*Hard Disc*” is spuriously correlated with “*label*”; however, “*label*” is not a background feature but rather part of the classification object. Spurious correlations also exist in **other data modalities**, e.g., text classification, leveraging the text generation capabilities of large language models (Brown et al., 2020). Other limitations of our work include evaluating **more generalization techniques** on Spawrious, including different robustness penalties (Liu et al., 2021; Blumberg et al., 2019; Krueger et al., 2021; Cha et al., 2021; Mahajan et al., 2021; Izmailov et al., 2022; Rame et al., 2022a), environment inference (Creager et al., 2021; Li et al., 2022; Sohoni et al., 2022; Huang et al., 2022), meta-learning (Zhang et al., 2020; Collins et al., 2020; Kaddour et al., 2020; Wang et al., 2021; Jiang et al., 2023), unsupervised domain adaptation (Ganin & Lempitsky, 2015; Long et al., 2016; Xu et al., 2021), dropout (LaBonte et al., 2022), flat minima (Cha et al., 2021; Kaddour et al., 2022a), weight averaging (Rame et al., 2022b; Wortsman et al., 2022; Kaddour, 2022), (counterfactual) data augmentation (Kaddour et al., 2022b; Goyal et al., 2021; Yao et al., 2022; Yin et al., 2023), fine-tuning of only specific layers (Kirichenko et al., 2022; Lee et al., 2023), diversity (Teney et al., 2022; Rame et al., 2022b), etc. Lastly, there is a possibility of **bias** creeping into the dataset via the generative model. Chuang et al. (2023) and others (Teo & Cheung, 2021; Zhao et al., 2018) have studied debiasing techniques for vision-language models, such as *Stable Diffusion v1*, and have moderate success in removing unexpected sources of spurious correlations.

7 CONCLUSION

We present Spawrious, an image classification benchmark with two types of spurious correlations, one-to-one (O2O) and many-to-many (M2M). We carefully design six dataset desiderata and instantiate them by leveraging recent advances in text-to-image and image captioning models. Next, we conduct experiments, and our findings indicate that even state-of-the-art group robustness techniques are insufficient in handling Spawrious, particularly in scenarios with Hard-splits where accuracy is below 73%. Our analysis of model errors revealed a dependence on irrelevant backgrounds, thus underscoring the difficulty of our dataset and highlighting the need for further investigations in this area. A more extensive discussion of limitations and future work can be found in Section 6.

REFERENCES

- 324
325
326 Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization.
327 *arXiv preprint arXiv:1907.02893*, 2019.
- 328 Martin Arjovsky, Kamalika Chaudhuri, and David Lopez-Paz. Throwing away data improves
329 worst-class error in imbalanced classification. *arXiv preprint arXiv:2205.11672*, 2022.
- 330
331 Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman
332 Vaughan. A theory of learning from different domains. *Machine learning*, 79:151–175, 2010.
- 333 Gilles Blanchard, Gyemin Lee, and Clayton Scott. Generalizing from several related classification
334 tasks to a new unlabeled sample. *Advances in neural information processing systems*, 24, 2011.
- 335
336 Stefano B. Blumberg, Marco Palombo, Can Son Khoo, Chantal M. W. Tax, Ryutaro Tanno, and
337 Daniel C. Alexander. Multi-stage prediction networks for data harmonization, 2019. URL
338 <https://arxiv.org/abs/1907.11629>.
- 339 Florian Bordes, Shashank Shekhar, Mark Ibrahim, Diane Bouchacourt, Pascal Vincent, and Ari S.
340 Morcos. Pug: Photorealistic and semantically controllable synthetic data for representation
341 learning, 2023.
- 342
343 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
344 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
345 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 346 Nicholas Carlini, Jamie Hayes, Milad Nasr, Matthew Jagielski, Vikash Sehwag, Florian Tramèr,
347 Borja Balle, Daphne Ippolito, and Eric Wallace. Extracting Training Data from Diffusion Models,
348 January 2023. URL <http://arxiv.org/abs/2301.13188>. arXiv:2301.13188 [cs].
- 349
350 Junbum Cha, Sanghyuk Chun, Kyungjae Lee, Han-Cheol Cho, Seunghyun Park, Yunsung Lee,
351 and Sungrae Park. Swad: Domain generalization by seeking flat minima. *Advances in Neural
352 Information Processing Systems*, 34:22405–22418, 2021.
- 353 Victor Chernozhukov, Denis Chetverikov, Mert Demirer, Esther Duflo, Christian Hansen, Whit-
354 ney Newey, and James Robins. Double/debiased machine learning for treatment and structural
355 parameters, 2018.
- 356
357 Mathieu Chevalley, Charlotte Bunne, Andreas Krause, and Stefan Bauer. Invariant causal mechanisms
358 through distribution matching. *arXiv preprint arXiv:2206.11646*, 2022.
- 359 Gordon Christie, Neil Fendley, James Wilson, and Ryan Mukherjee. Functional map of the world,
360 2017. URL <https://arxiv.org/abs/1711.07846>.
- 361
362 Ching-Yao Chuang, Varun Jampani, Yuanzhen Li, Antonio Torralba, and Stefanie Jegelka. Debiasing
363 vision-language models via biased prompts, 2023.
- 364 Liam Collins, Aryan Mokhtari, and Sanjay Shakkottai. Task-robust model-agnostic meta-learning.
365 *Advances in Neural Information Processing Systems*, 33:18860–18871, 2020.
- 366
367 Elliot Creager, Jörn-Henrik Jacobsen, and Richard Zemel. Environment inference for invariant
368 learning, 2021.
- 369 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
370 hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,
371 pp. 248–255. Ieee, 2009.
- 372
373 Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and
374 Zhifang Sui. A survey for in-context learning. *arXiv preprint arXiv:2301.00234*, 2022.
- 375 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
376 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An
377 image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint
arXiv:2010.11929*, 2020.

- 378 Qi Dou, Daniel Coelho de Castro, Konstantinos Kamnitsas, and Ben Glocker. Domain generalization
379 via model-agnostic learning of semantic features. *Advances in Neural Information Processing*
380 *Systems*, 32, 2019.
- 381 François-Guillaume Fernandez. Torchcam: class activation explorer. [https://github.com/
382 frgfm/torch-cam](https://github.com/frgfm/torch-cam), March 2020.
- 383 Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In
384 *International conference on machine learning*, pp. 1180–1189. PMLR, 2015.
- 385 Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and
386 Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias improves
387 accuracy and robustness. *arXiv preprint arXiv:1811.12231*, 2018a.
- 388 Robert Geirhos, Carlos RM Temme, Jonas Rauber, Heiko H Schütt, Matthias Bethge, and Felix A
389 Wichmann. Generalisation in humans and deep neural networks. *Advances in neural information
390 processing systems*, 31, 2018b.
- 391 Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias
392 Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature Machine
393 Intelligence*, 2(11):665–673, 2020.
- 394 Soumya Suvra Ghosal, Yifei Ming, and Yixuan Li. Are vision transformers robust to spurious
395 correlations?, 2022. URL <https://arxiv.org/abs/2203.09125>.
- 396 Sven Gowal, Sylvestre-Alvise Rebuffi, Olivia Wiles, Florian Stimberg, Dan Andrei Calian, and
397 Timothy A Mann. Improving robustness using generated data. *Advances in Neural Information
398 Processing Systems*, 34:4218–4233, 2021.
- 399 Ishaan Gulrajani and David Lopez-Paz. In search of lost domain generalization. In *International
400 Conference on Learning Representations*, 2021. URL [https://openreview.net/forum?
401 id=lQdXeXDwTlI](https://openreview.net/forum?id=lQdXeXDwTlI).
- 402 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
403 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
404 pp. 770–778, 2016.
- 405 Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common
406 corruptions and perturbations, 2019. URL <https://arxiv.org/abs/1903.12261>.
- 407 Katherine Hermann, Ting Chen, and Simon Kornblith. The origins and prevalence of texture
408 bias in convolutional neural networks. *Advances in Neural Information Processing Systems*, 33:
409 19000–19015, 2020.
- 410 Jeremy Howard. Imagenette: A smaller subset of 10 easily classified classes from imagenet, March
411 2019a. URL <https://github.com/fastai/imagenette>.
- 412 Jeremy Howard. Imagewoof: a subset of 10 classes from imagenet that aren’t so easy to classify,
413 March 2019b. URL <https://github.com/fastai/imagenette#imagewoof>.
- 414 Zeyi Huang, Haohan Wang, Dong Huang, Yong Jae Lee, and Eric P. Xing. The two dimensions of
415 worst-case training and the integrated effect for out-of-domain generalization, 2022.
- 416 Badr Youbi Idrissi, Martin Arjovsky, Mohammad Pezeshki, and David Lopez-Paz. Simple data
417 balancing achieves competitive worst-group-accuracy. In *Conference on Causal Learning and
418 Reasoning*, pp. 336–351. PMLR, 2022.
- 419 Pavel Izmailov, Polina Kirichenko, Nate Gruver, and Andrew Gordon Wilson. On feature learning in
420 the presence of spurious correlations. *arXiv preprint arXiv:2210.11369*, 2022.
- 421 Penghao Jiang, Ke Xin, Zifeng Wang, and Chunxi Li. Invariant meta learning for out-of-distribution
422 generalization. *arXiv preprint arXiv:2301.11779*, 2023.

- 432 Jean Kaddour. Stop wasting my time! saving days of imagenet and BERT training with latest
433 weight averaging. In *Has it Trained Yet? NeurIPS 2022 Workshop*, 2022. URL [https://](https://openreview.net/forum?id=00rABUHZuz)
434 openreview.net/forum?id=00rABUHZuz.
435
- 436 Jean Kaddour, Steindor Saemundsson, and Marc Deisenroth (he/him). Probabilistic Active Meta-
437 Learning. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin (eds.), *Ad-*
438 *vances in Neural Information Processing Systems*, volume 33, pp. 20813–20822. Curran As-
439 sociates, Inc., 2020. URL [https://proceedings.neurips.cc/paper/2020/file/](https://proceedings.neurips.cc/paper/2020/file/ef0d17b3bdb4ee2aa741ba28c7255c53-Paper.pdf)
440 [ef0d17b3bdb4ee2aa741ba28c7255c53-Paper.pdf](https://proceedings.neurips.cc/paper/2020/file/ef0d17b3bdb4ee2aa741ba28c7255c53-Paper.pdf).
- 441 Jean Kaddour, Yuchen Zhu, Qi Liu, Matt J Kusner, and Ricardo Silva. Causal effect inference for
442 structured treatments. *Advances in Neural Information Processing Systems*, 34:24841–24854,
443 2021.
- 444 Jean Kaddour, Linqing Liu, Ricardo Silva, and Matt Kusner. When do flat minima optimizers work?
445 In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in*
446 *Neural Information Processing Systems*, 2022a. URL [https://openreview.net/forum?](https://openreview.net/forum?id=vDeh2yxTvuh)
447 [id=vDeh2yxTvuh](https://openreview.net/forum?id=vDeh2yxTvuh).
448
- 449 Jean Kaddour, Aengus Lynch, Qi Liu, Matt J. Kusner, and Ricardo Silva. Causal machine learning:
450 A survey and open problems. *arXiv preprint arXiv:2206.15475*, 2022b. URL [https://arxiv.](https://arxiv.org/abs/2206.15475)
451 [org/abs/2206.15475](https://arxiv.org/abs/2206.15475).
- 452 Priyatham Kattakinda and Soheil Feizi. Focus: Familiar objects in common and uncommon settings,
453 2022.
454
- 455 Polina Kirichenko, Pavel Izmailov, and Andrew Gordon Wilson. Last layer re-training is sufficient
456 for robustness to spurious correlations. *arXiv preprint arXiv:2204.02937*, 2022.
457
- 458 Polina Kirichenko, Pavel Izmailov, and Andrew Gordon Wilson. Last layer re-training is sufficient
459 for robustness to spurious correlations, 2023.
- 460 Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Bal-
461 subramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, et al. Wilds: A
462 benchmark of in-the-wild distribution shifts. In *International Conference on Machine Learning*,
463 pp. 5637–5664. PMLR, 2021.
- 464 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep
465 convolutional neural networks. In F. Pereira, C. J. Burges, L. Bottou, and K. Q. Wein-
466 berger (eds.), *Advances in Neural Information Processing Systems*, volume 25. Curran Asso-
467 ciates, Inc., 2012. URL [https://proceedings.neurips.cc/paper_files/paper/](https://proceedings.neurips.cc/paper_files/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf)
468 [2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf).
469
- 470 David Krueger, Ethan Caballero, Joern-Henrik Jacobsen, Amy Zhang, Jonathan Binas, Dinghui
471 Zhang, Remi Le Priol, and Aaron Courville. Out-of-Distribution Generalization via Risk Extrapo-
472 lation (REx). *arXiv e-prints*, art. arXiv:2003.00688, March 2020. doi: 10.48550/arXiv.2003.00688.
- 473 David Krueger, Ethan Caballero, Joern-Henrik Jacobsen, Amy Zhang, Jonathan Binas, Dinghui
474 Zhang, Remi Le Priol, and Aaron Courville. Out-of-distribution generalization via risk extrapo-
475 lation (rex). In *International Conference on Machine Learning*, pp. 5815–5826. PMLR, 2021.
476
- 477 Sören R Künzel, Jasjeet S Sekhon, Peter J Bickel, and Bin Yu. Metalearners for estimating heteroge-
478 neous treatment effects using machine learning. *Proceedings of the national academy of sciences*,
479 116(10):4156–4165, 2019.
- 480 Tyler LaBonte, Vidya Muthukumar, and Abhishek Kumar. Dropout disagreement: A recipe for group
481 robustness with fewer annotations. In *NeurIPS 2022 Workshop on Distribution Shifts: Connecting*
482 *Methods and Applications*, 2022.
483
- 484 Yoonho Lee, Annie S. Chen, Fahim Tajwar, Ananya Kumar, Huaxiu Yao, Percy Liang, and Chelsea
485 Finn. Surgical Fine-Tuning Improves Adaptation to Distribution Shifts, March 2023. URL
<http://arxiv.org/abs/2210.11466>. arXiv:2210.11466 [cs].

- 486 Da Li, Jianshu Zhang, Yongxin Yang, Cong Liu, Yi-Zhe Song, and Timothy M Hospedales. Episodic
487 training for domain generalization. In *Proceedings of the IEEE/CVF International Conference on*
488 *Computer Vision*, pp. 1446–1455, 2019.
- 489 Haoliang Li, Sinno Jialin Pan, Shiqi Wang, and Alex C Kot. Domain generalization with adversarial
490 feature learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
491 pp. 5400–5409, 2018.
- 492 Zhiheng Li, Anthony Hoogs, and Chenliang Xu. Discover and mitigate unknown biases with
493 debiasing alternate networks, 2022.
- 494 Zhiheng Li, Ivan Evtimov, Albert Gordo, Caner Hazirbas, Tal Hassner, Cristian Canton Ferrer,
495 Chenliang Xu, and Mark Ibrahim. A whac-a-mole dilemma: Shortcuts come in multiples where
496 mitigating one amplifies others, 2023.
- 497 Weixin Liang and James Zou. Metashift: A dataset of datasets for evaluating contextual distribution
498 shifts and training conflicts. In *International Conference on Learning Representations*, 2022. URL
499 <https://openreview.net/forum?id=MTex8qKavoS>.
- 500 Evan Z Liu, Behzad Haghgoo, Annie S Chen, Aditi Raghunathan, Pang Wei Koh, Shiori Sagawa,
501 Percy Liang, and Chelsea Finn. Just train twice: Improving group robustness without training
502 group information. In *International Conference on Machine Learning*, pp. 6781–6792. PMLR,
503 2021.
- 504 Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In
505 *Proceedings of the IEEE international conference on computer vision*, pp. 3730–3738, 2015.
- 506 Mingsheng Long, Han Zhu, Jianmin Wang, and Michael I Jordan. Unsupervised domain adaptation
507 with residual transfer networks. *Advances in neural information processing systems*, 29, 2016.
- 508 Aengus Lynch, Jean Kaddour, and Ricardo Silva. Evaluating the impact of geometric and
509 statistical skews on out-of-distribution generalization performance. In *NeurIPS 2022 Work-*
510 *shop on Distribution Shifts: Connecting Methods and Applications*, 2022. URL <https://openreview.net/forum?id=wpT79coXAu>.
- 511 Divyat Mahajan, Shruti Tople, and Amit Sharma. Domain generalization using causal matching. In
512 *International Conference on Machine Learning*, pp. 7313–7324. PMLR, 2021.
- 513 Raghav Mehta, Vítor Albiero, Li Chen, Ivan Evtimov, Tamar Glaser, Zhiheng Li, and Tal Hass-
514 ner. You only need a good embeddings extractor to fix spurious correlations. *arXiv preprint*
515 *arXiv:2212.06254*, 2022.
- 516 Mazda Moayeri, Phillip Pope, Yogesh Balaji, and Soheil Feizi. A comprehensive study of image
517 classification model sensitivity to foregrounds, backgrounds, and visual attributes, 2022a.
- 518 Mazda Moayeri, Sahil Singla, and Soheil Feizi. Hard imagenet: Segmentations for objects with
519 strong spurious cues. *Advances in Neural Information Processing Systems*, 35:10068–10077,
520 2022b.
- 521 Krikamol Muandet, David Balduzzi, and Bernhard Schölkopf. Domain generalization via invariant
522 feature representation. In *International conference on machine learning*, pp. 10–18. PMLR, 2013.
- 523 Kevin P. Murphy. *Probabilistic Machine Learning: An introduction*. MIT Press, 2022. URL
524 probml.ai.
- 525 Vaishnavh Nagarajan, Anders Andreassen, and Behnam Neyshabur. Understanding the failure modes
526 of out-of-distribution generalization, 2020. URL <https://arxiv.org/abs/2010.15775>.
- 527 Yannic Neuhaus, Maximilian Augustin, Valentyn Boreiko, and Matthias Hein. Spurious features
528 everywhere – large-scale detection of harmful spurious features in imagenet, 2022. URL <https://arxiv.org/abs/2212.04871>.
- 529 Xinkun Nie and Stefan Wager. Quasi-oracle estimation of heterogeneous treatment effects. *Biometrika*,
530 108(2):299–319, 2021.

- 540 NLP Connect. vit-gpt2-image-captioning (revision 0e334c7), 2022. URL [https://](https://huggingface.co/nlpconnect/vit-gpt2-image-captioning)
541 huggingface.co/nlpconnect/vit-gpt2-image-captioning.
542
- 543 Daniel Omeiza, Skyler Speakman, Celia Cintas, and Komminist Weldermariam. Smooth grad-cam++:
544 An enhanced inference level visualization technique for deep convolutional neural network models,
545 2019.
- 546 Karl Pearson. Mathematical contributions to the theory of evolution.—on a form of spurious
547 correlation which may arise when indices are used in the measurement of organs. *Proceedings of*
548 *the royal society of london*, 60(359-367):489–498, 1897.
549
- 550 Jonas Peters, Peter Bühlmann, and Nicolai Meinshausen. Causal inference by using invariant
551 prediction: identification and confidence intervals. *Journal of the Royal Statistical Society Series*
552 *B: Statistical Methodology*, 78(5):947–1012, 2016.
- 553 Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of causal inference: foundations*
554 *and learning algorithms*. The MIT Press, 2017.
555
- 556 Alexandre Rame, Corentin Dancette, and Matthieu Cord. Fishr: Invariant gradient variances for out-
557 of-distribution generalization. In *International Conference on Machine Learning*, pp. 18347–18377.
558 PMLR, 2022a.
- 559 Alexandre Rame, Matthieu Kirchmeyer, Thibaud Rahier, Alain Rakotomamonjy, Patrick Gallinari,
560 and Matthieu Cord. Diverse weight averaging for out-of-distribution generalization. *arXiv preprint*
561 *arXiv:2205.09739*, 2022b.
- 562 Robin Rombach and Patrick Esser. License - a Hugging Face Space by CompVis, 2022. URL
563 <https://huggingface.co/spaces/CompVis/stable-diffusion-license>.
564
- 565 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
566 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Confer-*
567 *ence on Computer Vision and Pattern Recognition*, pp. 10684–10695, 2022.
- 568 Elan Rosenfeld, Pradeep Ravikumar, and Andrej Risteski. Domain-adjusted regression or: Erm
569 may already learn features sufficient for out-of-distribution generalization. *arXiv preprint*
570 *arXiv:2202.06856*, 2022.
571
- 572 Shiori Sagawa, Pang Wei Koh, Tatsunori B. Hashimoto, and Percy Liang. Distributionally robust
573 neural networks for group shifts: On the importance of regularization for worst-case generalization,
574 2019a. URL <https://arxiv.org/abs/1911.08731>.
- 575 Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust
576 neural networks for group shifts: On the importance of regularization for worst-case generalization.
577 *arXiv preprint arXiv:1911.08731*, 2019b.
578
- 579 Shiori Sagawa, Aditi Raghunathan, Pang Wei Koh, and Percy Liang. An investigation of why
580 overparameterization exacerbates spurious correlations. In *International Conference on Machine*
581 *Learning*, pp. 8346–8356. PMLR, 2020.
- 582 Pamela Samuelson. Generative AI meets copyright. *Science*, 381(6654):158–161, July 2023. doi:
583 10.1126/science.adi0656. URL [https://www.science.org/doi/10.1126/science.](https://www.science.org/doi/10.1126/science.adi0656)
584 [adi0656](https://www.science.org/doi/10.1126/science.adi0656). Publisher: American Association for the Advancement of Science.
585
- 586 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi
587 Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski,
588 Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev.
589 Laion-5b: An open large-scale dataset for training next generation image-text models, 2022. URL
590 <https://arxiv.org/abs/2210.08402>.
- 591 Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh,
592 and Dhruv Batra. Grad-CAM: Visual explanations from deep networks via gradient-based lo-
593 calization. *International Journal of Computer Vision*, 128(2):336–359, oct 2019. doi: 10.1007/
s11263-019-01228-7. URL <https://doi.org/10.1007%2Fs11263-019-01228-7>.

- 594 Yuge Shi, Jeffrey Seely, Philip H. S. Torr, N. Siddharth, Awni Hannun, Nicolas Usunier, and Gabriel
595 Synnaeve. Gradient matching for domain generalization. *arXiv preprint arXiv:2104.09937*, 2021.
596
- 597 Herbert A Simon. Spurious correlation: A causal interpretation. *Journal of the American statistical*
598 *Association*, 49(267):467–479, 1954.
- 599 Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks:
600 Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*, 2013.
601
- 602 Sahil Singla and Soheil Feizi. Salient imagenet: How to discover spurious features in deep learning?,
603 2022.
604
- 605 Nimit S. Sohoni, Jared A. Dunnmon, Geoffrey Angus, Albert Gu, and Christopher Ré. No subclass
606 left behind: Fine-grained robustness in coarse-grained classification problems, 2022.
- 607 Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Under-
608 standing and Mitigating Copying in Diffusion Models, May 2023. URL [http://arxiv.org/](http://arxiv.org/abs/2305.20086)
609 [abs/2305.20086](http://arxiv.org/abs/2305.20086). arXiv:2305.20086 [cs].
610
- 611 Baochen Sun and Kate Saenko. Deep coral: Correlation alignment for deep domain adaptation. In
612 *Computer Vision—ECCV 2016 Workshops: Amsterdam, The Netherlands, October 8–10 and 15–16,*
613 *2016, Proceedings, Part III 14*, pp. 443–450. Springer, 2016.
- 614 Damien Teney, Ehsan Abbasnejad, Simon Lucey, and Anton Van den Hengel. Evading the simplicity
615 bias: Training a diverse set of models discovers solutions with superior ood generalization. In
616 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.
617 16761–16772, 2022.
618
- 619 Christopher T. H Teo and Ngai-Man Cheung. Measuring fairness in generative models, 2021.
- 620 Vladimir Vapnik. Principles of risk minimization for learning theory. *Advances in neural information*
621 *processing systems*, 4, 1991.
622
- 623 Joshua Vendrow, Saachi Jain, Logan Engstrom, and Aleksander Madry. Dataset interfaces: Diagnos-
624 ing model failures using controllable counterfactual generation. *arXiv preprint arXiv:2302.07865*,
625 2023.
626
- 627 Zhenyi Wang, Tiehang Duan, Le Fang, Qiuling Suo, and Mingchen Gao. Meta learning on a sequence
628 of imbalanced domains with difficulty awareness. In *Proceedings of the IEEE/CVF International*
629 *Conference on Computer Vision*, pp. 8947–8957, 2021.
- 630 Olivia Wiles, Sven Gowal, Florian Stimberg, Sylvestre Alvisé-Rebuffi, Ira Ktena, Krishnamurthy
631 Dvijotham, and Taylan Cemgil. A fine-grained analysis on distribution shift. *arXiv preprint*
632 *arXiv:2110.11328*, 2021.
633
- 634 Olivia Wiles, Isabela Albuquerque, and Sven Gowal. Discovering bugs in vision models using
635 off-the-shelf image generation and captioning. *arXiv preprint arXiv:2208.08831*, 2022.
- 636 Mitchell Wortsman, Gabriel Ilharco, Samir Ya Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes,
637 Ari S Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, et al. Model
638 soups: averaging weights of multiple fine-tuned models improves accuracy without increasing
639 inference time. In *International Conference on Machine Learning*, pp. 23965–23998. PMLR,
640 2022.
- 641 Kai Xiao, Logan Engstrom, Andrew Ilyas, and Aleksander Madry. Noise or signal: The role of image
642 backgrounds in object recognition. *arXiv preprint arXiv:2006.09994*, 2020.
643
- 644 Minghao Xu, Jian Zhang, Bingbing Ni, Teng Li, Chengjie Wang, Qi Tian, and Wenjun Zhang.
645 Adversarial domain adaptation with domain mixup, 2019.
646
- 647 Tongkun Xu, Weihua Chen, Pichao Wang, Fan Wang, Hao Li, and Rong Jin. Cdtrans: Cross-domain
transformer for unsupervised domain adaptation. *arXiv preprint arXiv:2109.06165*, 2021.

648 Huaxiu Yao, Yu Wang, Sai Li, Linjun Zhang, Weixin Liang, James Zou, and Chelsea Finn. Im-
649 proving Out-of-Distribution Robustness via Selective Augmentation. In *Proceedings of the 39th*
650 *International Conference on Machine Learning*, pp. 25407–25437. PMLR, June 2022. URL
651 <https://proceedings.mlr.press/v162/yao22b.html>. ISSN: 2640-3498.
652

653 Nanyang Ye, Kaican Li, Haoyue Bai, Runpeng Yu, Lanqing Hong, Fengwei Zhou, Zhenguo Li,
654 and Jun Zhu. Ood-bench: Quantifying and understanding two dimensions of out-of-distribution
655 generalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
656 *Recognition*, pp. 7947–7958, 2022.

657 Yuwei Yin, Jean Kaddour, Xiang Zhang, Yixin Nie, Zhenguang Liu, Lingpeng Kong, and Qi Liu.
658 Ttida: Controllable generative data augmentation via text-to-text and text-to-image models, 2023.
659

660 Marvin Zhang, Henrik Marklund, Abhishek Gupta, Sergey Levine, and Chelsea Finn. Adaptive risk
661 minimization: A meta-learning approach for tackling group shift. *arXiv preprint arXiv:2007.02931*,
662 8:9, 2020.

663 Xingxuan Zhang, Yue He, Tan Wang, Jiabin Qi, Han Yu, Zimu Wang, Jie Peng, Renzhe Xu, Zheyang
664 Shen, Yulei Niu, et al. Nico challenge: Out-of-distribution generalization for image recognition
665 challenges. In *European Conference on Computer Vision*, pp. 433–450. Springer, 2023.

666 Shengjia Zhao, Hongyu Ren, Arianna Yuan, Jiaming Song, Noah Goodman, and Stefano Ermon.
667 Bias and generalization in deep generative models: An empirical study, 2018.
668

669 Evan Zheran Liu, Behzad Haghgoo, Annie S. Chen, Aditi Raghunathan, Pang Wei Koh, Shiori
670 Sagawa, Percy Liang, and Chelsea Finn. Just Train Twice: Improving Group Robustness without
671 Training Group Information. *arXiv e-prints*, art. arXiv:2107.09044, July 2021. doi: 10.48550/
672 arXiv.2107.09044.

673 Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep
674 features for discriminative localization, 2015.
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

A EXISTING BENCHMARKS

We summarize the differences between Spawrious and related benchmarks in Table 3. DomainBed (Gulrajani & Lopez-Paz, 2021) is a benchmark suite consisting of seven previously published datasets focused on domain generalization (DG), not on spurious correlations (excluding CMNIST, which we discuss separately). After careful hyper-parameter tuning, the authors find that ERM, not specifically designed for DG settings, as well as DG-specific methods, perform all about the same on average. They conjecture that these datasets may comprise an ill-posed challenge. For example, they raise the question of whether DG from a photo-realistic training environment to a cartoon test environment is even possible. In contrast, we follow the same rigorous hyper-parameter tuning procedure by (Gulrajani & Lopez-Paz, 2021) and observe stark differences among methods on Spawrious in Section 4, with ERM being the worst and 10.68% points worse than the best method on average.

Like DomainBed, OoD-Bench (Ye et al., 2022) combines previously published datasets with the added contribution of characterizing them as a combination of diversity shift and style shift, allowing the evaluation of algorithms on a more comprehensive range of shifts. Methods that handle both shifts, like (Huang et al., 2022), will consistently beat ERM. By testing on unseen backgrounds-foreground combinations while having correlated backgrounds, we can address the two types of shifts they describe, while most datasets only address one type of shift. WILDS (Koh et al., 2021), NICO (Zhang et al., 2023), FOCUS (Kattakinda & Feizi, 2022), MetaShift (Liang & Zou, 2022) collect in-the-wild data and group data points with environment labels. However, these benchmarks do not induce *explicit* spurious correlations between environments and labels. For example, WILDS-FMOW (Koh et al., 2021; Christie et al., 2017) possesses a label shift between non-African and African regions; yet, the test images pose a domain generalization (DG) challenge (test images were taken several years later than training images) instead of reverting the spurious correlations observed in the training data. Waterbirds (Sagawa et al., 2019a), and CelebA hair color (Liu et al., 2015; Izmailov et al., 2022) are binary classification datasets including spurious correlations but without unseen test domains (DG). Further, Idrissi et al. (2022) illustrates that a simple class-balancing strategy alleviates most of their difficulty, while Spawrious is class-balanced from the beginning. ColorMNIST (Arjovsky et al., 2019) includes spurious correlations and poses a DG problem. However, it is based on MNIST and, therefore, over-simplistic, i.e., it does not reflect real-world spurious correlations involving complex background features, such as the ones found in ImageNet (Singla & Feizi, 2022; Neuhaus et al., 2022). Hard ImageNet (Moayeri et al., 2022b) is a benchmark created by collecting images in ImageNet that contain spurious features, however, they do not satisfy our desiderata of multiple training environments and multiple difficulty levels Section 2. Like us, Li et al. (2023) create two synthetic datasets, UrbanCars and ImageNet-W, to test for spurious feature reliance, but these datasets do not satisfy our desiderata of photorealism and high-fidelity backgrounds Section 2. PUG (Bordes et al., 2023) synthetically generate a dataset of unfamiliar object-location images, but they do not create a benchmark that introduces *explicit* spurious correlations between environment and labels. None of the above benchmarks include explicit training and test environments for M2M-SCs.

B MORE RELATED WORK

Causal Inference The theory of causation provides another perspective on the sources and possible mitigations of spurious correlations (Peters et al., 2016; 2017; Kaddour et al., 2022b). Namely, we can formalize environment-specific data as samples from different interventional distributions, which keep the influence of variables not affected by the corresponding interventions invariant. This perspective has motivated several invariance-learning methods that make causal assumptions on the data-generating process (Arjovsky et al., 2019; Kaddour et al., 2022b). The field of treatment effect estimation also deals with mitigating spurious correlations from observational data (Chernozhukov et al., 2018; Künzel et al., 2019; Kaddour et al., 2021; Nie & Wager, 2021).

Dataset	DG	O2O-SC	M2M-SC	Synthetic	Dataset Size
CelebA-Hair Color Liu et al. (2015)	X	✓	X	X	162770
Waterbirds Sagawa et al. (2019a)	X	✓	X	✓	4795
CMNIST Arjovsky et al. (2019)	✓	✓	X	✓	60000
DomainBed* Gulrajani & Lopez-Paz (2021)	✓	X	X	X	-
WILDS Koh et al. (2021)	✓	X	X	X	-
NICO Zhang et al. (2023)	✓	X	X	X	25000
MetaShift Liang & Zou (2022)	✓	X	X	X	12868
Spawrious	✓	✓	✓	✓	152000

Table 3: **Differences between Spawrious and other benchmarks**, according to whether they pose a Domain Generalization (DG), One-To-One-and/or Many-To-Many Spurious Correlations challenge.

Test-time domain adaptation with labels involves either fine-tuning a model [Rosenfeld et al. \(2022\)](#); [Izmailov et al. \(2022\)](#); [Kirichenko et al. \(2023\)](#) or in-context learning [Dong et al. \(2022\)](#) to leverage a small amount of labeled test-domain examples.

Miscellaneous [Nagarajan et al. \(2020\)](#) analyze two different kinds of spurious correlations: *geometric* and *statistical* skew. Geometric skew occurs when there is an imbalance between groups of types of data points (i.e., data points from different environments) and leads to misclassification when the balance of groups changes. This understanding has motivated simply removing data points from the training data to balance between groups of data points ([Arjovsky et al., 2022](#)). In contrast, we study two particular types of SCs, which persist in degenerating generalization performance despite perfect balances of classes among groups. Further, [Ye et al. \(2022\)](#) provide a two-dimensional decomposition of OOD difficulty into correlation and diversity shifts between the training and test set. The challenges in our work span both of these dimensions, because the test environment contains unseen background-foreground combinations, a diversity shift, and the background is spuriously correlated with the foreground in the training data, a correlation shift.

C EXPERIMENTAL SETUP

Methods The field of worst-group-accuracy optimization is thriving with a plethora of proposed methods, making it impractical to compare all available methods. We choose the following six popular methods and their `DomainBed` implementation ([Gulrajani & Lopez-Paz, 2021](#)). **ERM** ([Vapnik, 1991](#)) refers to the canonical, average-accuracy-optimization procedure, where we treat all groups identically and ignore group labels, not targeting to improve the worst group performance. **CORAL** ([Sun & Saenko, 2016](#)) penalizes differences in the first and second moment of the feature distributions of each group. **IRM** ([Arjovsky et al., 2019](#)) is a causality-inspired ([Kaddour et al., 2022b](#)) invariance-learning method, which penalizes feature distributions that have different optimal linear classifiers for each group. **CausIRL** ([Chevalley et al., 2022](#)) is another causally-motivated algorithm for learning invariances, whose penalty considers only one distance between mixtures of latent features coming from different domains. **GroupDRO** ([Sagawa et al., 2019a](#)) uses Group-Distributional Robust Optimization to explicitly minimize the worst group loss instead of the average loss. **MMD-AAE** ([Li et al., 2018](#)) penalizes distances between feature distributions of groups via the maximum mean discrepancy (MMD) and learning an adversarial auto-encoder (AAE). **JTT** ([Zheran Liu et al., 2021](#)) runs ERM for a certain number of epochs, stops, then runs classifications on all the training samples; then the misclassifications are up-weighted in the loss, and training continues. **W2D** ([Huang et al., 2022](#)) upweights datapoints in the loss that have either high *feature loss* or *sample loss*. **VREx** ([Krueger et al., 2020](#)) penalizes variance between the environment-specific training losses. **Fish** ([Shi et al., 2021](#)) rewards large inner products between environment-specific training gradients. **Mixup** ([Xu et al., 2019](#)) linearly interpolates between two images’ pixel values, and has been implemented with random shuffle (randomly mix images across environments and labels) and **LISA** ([Yao et al., 2022](#)) (alternate between mixing across environments for the same label, or across labels for the same environment).

Hyper-parameter tuning We follow the hyper-parameter tuning process used in `DomainBed` ([Gulrajani & Lopez-Paz, 2021](#)) with a minor modification. We keep the dropout rate (0.1) and the batch size fixed (128 for ResNets and 64 for ViTs) because we found them to have only a very marginal impact on the performance. We tune the learning rate and weight decay on ERM with a random search of 20 random trials. For all other methods, we further tune their method-specific hyper-parameters with a search of 10 random trials. We perform model selection based on the training domain validation accuracy of a subset of the training data. We reuse the hyper-parameters found for Spawrious-`{O2O}`-`{Easy}` and Spawrious-`{M2M}`-`{Hard}` on Spawrious-`{O2O}`-`{Medium, Hard}` and Spawrious-`{M2M}`-`{Easy, Medium}`, respectively. We also initially explored the ViT ([Dosovitskiy et al., 2020](#)) architecture, with results shown in Appendix F. Due to its poor performance, we chose to focus on ResNet50 results.

Evaluation We evaluate the classifiers on a test environment where the SCs present during training change, as described in Table 1. For O2O, multiple ways exist to choose a test data combination; we evaluate one of them as selected using a random search process. In M2M, because there are only two class groups and two background groups, we only need to swap them as seen in Figure 1b.

D ETHICAL CONCERNS

D.1 BIASES

We first acknowledged that generative models can inherit biases from their training data, including those related to dog breed representation and dog breed characteristics. We utilized various measures to mitigate these biases:

- *Dog Breed Representation*: By design, we ensured that the breeds in our dataset are balanced, avoiding underrepresentation or overrepresentation of any particular breed.
- *Dog Breed Characteristics*: We examined the characteristics associated with each breed and verified that our model does not exaggerate or stereotype them.

Further, we employed quality control measures, as described in Section 4.1, to guarantee that images are realistic and high-quality, regardless of breed. We manually reviewed the generated images to ensure they were free from harmful associations and stereotypes.

D.2 COPYRIGHT CONSIDERATIONS

We purposefully decided to use StableDiffusion, which offers a permissive license that allows for commercial and non-commercial usage. See more info in (Rombach & Esser, 2022).

Further, we are aware of possible copyright and fair use offenses, which are still debated. To our knowledge, under US law, fair uses of in-copyright works do not infringe copyrights Samuelson (2023). Courts consider four factors when assessing fair use defenses: (1) the purpose of the challenged use, (2) the nature of the copyrighted works, (3) the amount and substantiality of the taking, and (4) the effect of the challenged use on the market for or value of the copyrighted work, which we address as follows:

1. *Purpose and character*: Academic research is nonprofit and educational.
2. *Nature of the work*: Academic research often involves factual or informational works.
3. *Amount and substantiality*: We use generated images, which are likely to include only small portions if any of copyrighted works (Carlini et al., 2023; Somepalli et al., 2023).
4. *Effect on the market*: Academic research is unlikely to harm the market for the original work.

E EFFECT OF IMAGENET PRE-TRAINING

Method	One-To-One SC			Many-To-Many SC			Average
	Easy	Medium	Hard	Easy	Medium	Hard	
ERM	45.75%±1.26	46.86%±1.10	41.85% ±0.56	57.67%±2.55	30.03%±0.28	30.05%±1.34	42.04%
GroupDRO	46.50% ±0.91	46.52%±0.95	39.80%±1.66	60.82%±0.58	31.72%±0.35	31.62% ±1.72	42.83%
MMD-AAE	44.09%±1.80	46.87% ±1.46	39.67%±0.84	61.24% ±0.93	32.10% ±0.47	30.77%±1.58	42.46%
ERM	77.49%±0.05	76.60% ±0.02	71.32%±0.09	83.80% ±0.01	53.05%±0.03	58.70%±0.04	70.16%
GroupDRO	80.58% ±0.74	75.96%±2.18	76.99% ±2.60	79.96%±2.79	61.01% ±4.64	60.86% ±1.71	72.56%
MMD-AAE	78.81%±0.02	75.33%±0.03	72.66%±0.01	80.55%±0.02	59.43%±0.04	54.39%±0.05	70.20%

Table 4: **Impact of ImageNet pretraining**: ResNet-50 without ImageNet pretraining (top) vs ResNet-50 with ImageNet pretraining (bottom) results

We have included ImageNet pretraining for all of our main body results in Table 3, as has been done for results comparisons on Waterbirds (Sagawa et al., 2019a) and CelebA (Liu et al., 2015) and has become standard practice for image classification (Krizhevsky et al., 2012). However, we also measure the performance of a ResNet50 trained just on the Spawrious challenges and report our results in Table 4. We find that pretraining makes a consistently positive impact on the performance of the classifiers, with a 28.12% point difference between the ERM performances.

864
865
866
867
868
869
870
871
872
873
874
875



876 Figure 4: **ERM misclassifications due to spurious correlations.** The shown test images correspond
877 to the class "Bulldog" with spurious backgrounds "Mountains" in the O2O-Hard (left) and "Snow"
878 in the M2M-Hard (right) challenge.

879
880
881

F EFFECT OF MODEL ARCHITECTURE

882
883
884
885
886
887
888
889
890

Method	One-To-One SC			Many-To-Many SC			Average
	Easy	Medium	Hard	Easy	Medium	Hard	
ERM	36.28%±1.17	32.78%±2.55	30.2%±0.83	55.56%±0.75	32.78%±2.55	30.20%±0.83	40.44%
GroupDRO	41.14%±1.62	51.43%±0.53	40.21%±1.76	53.79%±1.35	30.79%±1.75	25.45%±1.15	40.47%
MMD-AAE	40.64%±3.11	53.36%±0.95	38.54%±1.92	58.42%±1.77	24.75%±0.59	28.91%±2.68	40.77%
ERM	77.49%±0.05	76.60%±0.02	71.32%±0.09	83.80%±0.01	53.05%±0.03	58.70%±0.04	70.16%
GroupDRO	80.58%±0.74	75.96%±2.18	76.99%±2.60	79.96%±2.79	61.01%±4.64	60.86%±1.71	72.56%
MMD-AAE	78.81%±0.02	75.33%±0.03	72.66%±0.01	80.55%±0.02	59.43%±0.04	54.39%±0.05	70.20%

891 Table 5: **Impact of ViT-B instead of ResNet-50:** ViT-B pretrained on ImageNet (top) vs ResNet-50
892 pretrained on ImageNet (bottom) results

893
894
895
896
897
898
899
900
901

We experiment with the ViT-B/16 (Dosovitskiy et al., 2020), following (Izmailov et al., 2022; Mehta et al., 2022). Based on Table 5, we make the following observations: The ViT backbone architecture worsens the performance for both MMD-AAE and ERM, underperforming the ResNet50. The best results for ERM were obtained with ResNet50, which performs 29.72% points better than the best ViT. In the debate on whether ViTs (Dosovitskiy et al., 2020) are generally more robust to SCs (Ghosal et al., 2022) than CNNs or not (Izmailov et al., 2022; Mehta et al., 2022), our results side with the latter. We observe that a ViT-B/16 pretrained on ImageNet22k had worse test accuracies than the ResNet architecture.

902
903
904

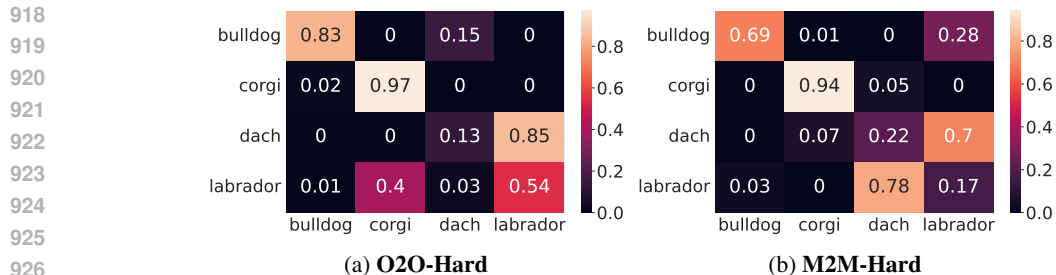
G MISCLASSIFICATIONS ANALYSIS

905
906
907
908
909
910
911
912
913
914
915
916
917

In Section 4, we learned that ERM performs particularly poorly on both hard challenges. Now, we want to investigate why by examining some of the misclassifications. For example, we observe in Figure 4 that on the test set, the class "Bulldog" is misclassified as the classes whose most common training set background is the same as "Bulldog"'s test backgrounds.

Note that for all classes and in all data groups, both training and test environments, the number of data points per class is always balanced; rendering methods like *Subsampling large classes* (Idrissi et al., 2022), which achieve state-of-the-art performance on other SC benchmarks, inapplicable. Hence, we conjecture that despite balanced classes, the model heavily relies on the spurious features of the "Mountains" and "Snow" backgrounds.

We further corroborate that claim by examining the model's confusion matrix in Figure 5. For example, Figure 5a shows the highest non-diagonal value for actual "Dachshund" images being wrongly classified as "Labrador". We conjecture the reason being that in O2O-Hard, the background of "Dachshund" in the test set is "Snow", which is the most common background of the training



928 **Figure 5: Confusion matrices for ERM models.** *X*-axis: predictions; *Y*-axis: true labels.

929

930

931 images of “*Labrador*”, as shown in Table 1. We examine the features learned by the ERM model

932 using saliency maps in Appendix G.

933 Saliency maps (Simonyan et al., 2013; Zhou et al., 2015; Selvaraju et al., 2019; Omeiza et al., 2019)

934 are a method for investigating the input features that most positively affect a model’s particular

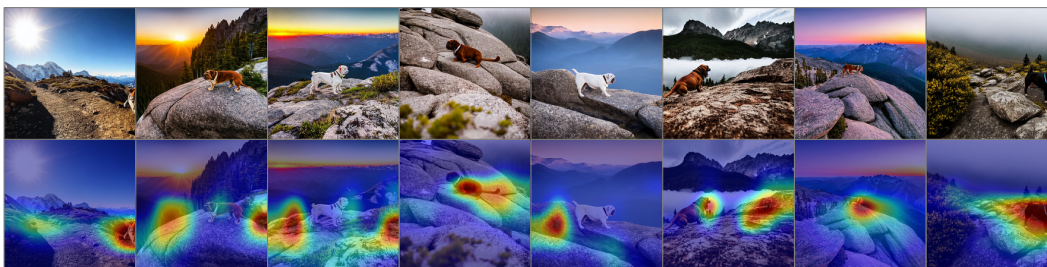
935 classification. We applied the Smooth Grad-CAM++ saliency map method (Omeiza et al., 2019;

936 Fernandez, 2020) to the misclassified images from an ERM model in the test domains of the O2O-

937 Hard and M2M-Hard challenges. The saliency maps we obtained in Figure 6 and Figure 7 suggest

938 that the ERM model was sensitive to (spurious) background features, although seemingly more in the

939 O2O challenge than the M2M challenge.



949 **Figure 6: O2O-Hard saliency maps:** all images were misclassifications of *Bulldog* as *Dachshund*



962 **Figure 7: M2M-Hard saliency maps:** all images were misclassifications of *Bulldog* as *Labrador*

963

964

965 Next, we compare qualitatively the difference in saliency maps between the Mixup and ERM

966 optimization methods, which can be seen in Figure 8. While the exact saliency pattern differs between

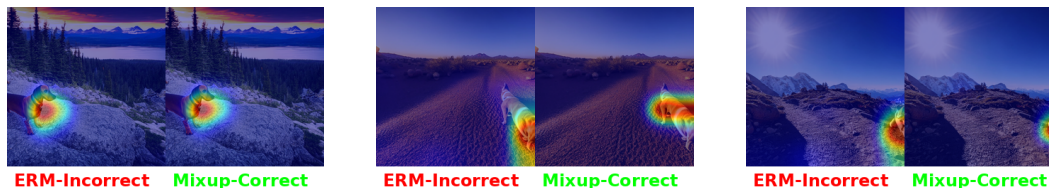
967 the two methods, they ultimately seem to be attending to the same image features.

968 H FAILURE ANALYSIS OF THE GENERATION PIPELINE

969

970

971 We conduct a failure analysis in two ways: manual and automatic. In our manual visual examination, we inspected large samples of the generated images via human annotators (the authors). Our



979 **Figure 8: Saliency comparisons between Mixup and ERM**

980
981
982 automated failure analysis pipeline is described in Section 4.2. For example, to test the quality of
983 a prompt, we only accept it under two conditions: at least 95 images out of 100 look realistic and
984 fit the prompt. Second, all remaining images must only be unfit because of the absence of a dog in
985 the image. Identifying a dog in an image is a relatively easy task for the image captioning model.
986 We confirmed by evaluating on the unfit images and assessing that they all get flagged by the image
987 captioning model (the caption does not contain the word dog).
988

989 I CLEANLINESS ANALYSIS OF THE DATASET



1004
1005 **Figure 9: Volunteers decided on prompt-image alignment for 224x224 images:** We asked 10
1006 volunteers to scan images such as the three shown above and return a score for the number of correctly
1007 aligned images
1008

1009 We have checked the accuracy of prompt-image alignment of images such as those in Figure 9 from
1010 a random sample of our dataset using human annotators (10 volunteers). We collected a random
1011 sample of 480 images from our dataset, appended with the intended caption for the image, and then
1012 partitioned this dataset into 10 folders. We asked 10 volunteers to scan the images and return a score
1013 for the number of correctly aligned images. Our scores were: 48, 46, 46, 46, 47, 47, 46, 46, 47, 48;
1014 resulting in an average of $46.7/48 = 97.2\%$.
1015

1016 J DISCUSSION OF M2M VS O2O

1017
1018 In order to understand how the M2M challenge leads to poor generalisation performance, consider
1019 the following situation, where the classifier achieves low loss in training by simulating a decision
1020 tree within the network, as depicted in Figure 2b of the submission. The model first represents the
1021 background, and then decides which group of dogs the image could be representing conditioned on
1022 the background. Within this setting, the spurious feature dependence arises at the beginning of the
1023 decision tree. In the test data, this decision tree fails to work because the background group is wholly
1024 unpredictable of the class groups. As seen in Figure 2d, the blue background group (s3, s4) is a feature
1025 used by the model to decide between classes (c3, c4), when in fact the model should be deciding
between (c1, c2).

1026 REFERENCES

- 1027
1028 Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization.
1029 *arXiv preprint arXiv:1907.02893*, 2019.
- 1030 Martin Arjovsky, Kamalika Chaudhuri, and David Lopez-Paz. Throwing away data improves
1031 worst-class error in imbalanced classification. *arXiv preprint arXiv:2205.11672*, 2022.
- 1032
1033 Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman
1034 Vaughan. A theory of learning from different domains. *Machine learning*, 79:151–175, 2010.
- 1035 Gilles Blanchard, Gyemin Lee, and Clayton Scott. Generalizing from several related classification
1036 tasks to a new unlabeled sample. *Advances in neural information processing systems*, 24, 2011.
- 1037
1038 Stefano B. Blumberg, Marco Palombo, Can Son Khoo, Chantal M. W. Tax, Ryutaro Tanno, and
1039 Daniel C. Alexander. Multi-stage prediction networks for data harmonization, 2019. URL
1040 <https://arxiv.org/abs/1907.11629>.
- 1041 Florian Bordes, Shashank Shekhar, Mark Ibrahim, Diane Bouchacourt, Pascal Vincent, and Ari S.
1042 Morcos. Pug: Photorealistic and semantically controllable synthetic data for representation
1043 learning, 2023.
- 1044
1045 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
1046 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
1047 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 1048 Nicholas Carlini, Jamie Hayes, Milad Nasr, Matthew Jagielski, Vikash Sehwal, Florian Tramèr,
1049 Borja Balle, Daphne Ippolito, and Eric Wallace. Extracting Training Data from Diffusion Models,
1050 January 2023. URL <http://arxiv.org/abs/2301.13188>. arXiv:2301.13188 [cs].
- 1051
1052 Junbum Cha, Sanghyuk Chun, Kyungjae Lee, Han-Cheol Cho, Seunghyun Park, Yunsung Lee,
1053 and Sungrae Park. Swad: Domain generalization by seeking flat minima. *Advances in Neural
1054 Information Processing Systems*, 34:22405–22418, 2021.
- 1055 Victor Chernozhukov, Denis Chetverikov, Mert Demirer, Esther Duflo, Christian Hansen, Whit-
1056 ney Newey, and James Robins. Double/debiased machine learning for treatment and structural
1057 parameters, 2018.
- 1058
1059 Mathieu Chevalley, Charlotte Bunne, Andreas Krause, and Stefan Bauer. Invariant causal mechanisms
1060 through distribution matching. *arXiv preprint arXiv:2206.11646*, 2022.
- 1061
1062 Gordon Christie, Neil Fendley, James Wilson, and Ryan Mukherjee. Functional map of the world,
2017. URL <https://arxiv.org/abs/1711.07846>.
- 1063
1064 Ching-Yao Chuang, Varun Jampani, Yuanzhen Li, Antonio Torralba, and Stefanie Jegelka. Debiasing
1065 vision-language models via biased prompts, 2023.
- 1066
1067 Liam Collins, Aryan Mokhtari, and Sanjay Shakkottai. Task-robust model-agnostic meta-learning.
Advances in Neural Information Processing Systems, 33:18860–18871, 2020.
- 1068
1069 Elliot Creager, Jörn-Henrik Jacobsen, and Richard Zemel. Environment inference for invariant
1070 learning, 2021.
- 1071
1072 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
1073 hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,
pp. 248–255. Ieee, 2009.
- 1074
1075 Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and
1076 Zhifang Sui. A survey for in-context learning. *arXiv preprint arXiv:2301.00234*, 2022.
- 1077
1078 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
1079 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An
image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint
arXiv:2010.11929*, 2020.

- 1080 Qi Dou, Daniel Coelho de Castro, Konstantinos Kamnitsas, and Ben Glocker. Domain generalization
1081 via model-agnostic learning of semantic features. *Advances in Neural Information Processing*
1082 *Systems*, 32, 2019.
- 1083 François-Guillaume Fernandez. Torchcam: class activation explorer. [https://github.com/
1084 frgfm/torch-cam](https://github.com/frgfm/torch-cam), March 2020.
- 1085 Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In
1086 *International conference on machine learning*, pp. 1180–1189. PMLR, 2015.
- 1087 Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and
1088 Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias improves
1089 accuracy and robustness. *arXiv preprint arXiv:1811.12231*, 2018a.
- 1090 Robert Geirhos, Carlos RM Temme, Jonas Rauber, Heiko H Schütt, Matthias Bethge, and Felix A
1091 Wichmann. Generalisation in humans and deep neural networks. *Advances in neural information
1092 processing systems*, 31, 2018b.
- 1093 Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias
1094 Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature Machine
1095 Intelligence*, 2(11):665–673, 2020.
- 1096 Soumya Suvra Ghosal, Yifei Ming, and Yixuan Li. Are vision transformers robust to spurious
1097 correlations?, 2022. URL <https://arxiv.org/abs/2203.09125>.
- 1098 Sven Gowal, Sylvestre-Alvise Rebuffi, Olivia Wiles, Florian Stimberg, Dan Andrei Calian, and
1099 Timothy A Mann. Improving robustness using generated data. *Advances in Neural Information
1100 Processing Systems*, 34:4218–4233, 2021.
- 1101 Ishaan Gulrajani and David Lopez-Paz. In search of lost domain generalization. In *International
1102 Conference on Learning Representations*, 2021. URL [https://openreview.net/forum?
1103 id=lQdXeXDwTl](https://openreview.net/forum?id=lQdXeXDwTl).
- 1104 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
1105 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
1106 pp. 770–778, 2016.
- 1107 Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common
1108 corruptions and perturbations, 2019. URL <https://arxiv.org/abs/1903.12261>.
- 1109 Katherine Hermann, Ting Chen, and Simon Kornblith. The origins and prevalence of texture
1110 bias in convolutional neural networks. *Advances in Neural Information Processing Systems*, 33:
1111 19000–19015, 2020.
- 1112 Jeremy Howard. Imagenette: A smaller subset of 10 easily classified classes from imagenet, March
1113 2019a. URL <https://github.com/fastai/imagenette>.
- 1114 Jeremy Howard. Imagewoof: a subset of 10 classes from imagenet that aren’t so easy to classify,
1115 March 2019b. URL <https://github.com/fastai/imagenette#imagewoof>.
- 1116 Zeyi Huang, Haohan Wang, Dong Huang, Yong Jae Lee, and Eric P. Xing. The two dimensions of
1117 worst-case training and the integrated effect for out-of-domain generalization, 2022.
- 1118 Badr Youbi Idrissi, Martin Arjovsky, Mohammad Pezeshki, and David Lopez-Paz. Simple data
1119 balancing achieves competitive worst-group-accuracy. In *Conference on Causal Learning and
1120 Reasoning*, pp. 336–351. PMLR, 2022.
- 1121 Pavel Izmailov, Polina Kirichenko, Nate Gruver, and Andrew Gordon Wilson. On feature learning in
1122 the presence of spurious correlations. *arXiv preprint arXiv:2210.11369*, 2022.
- 1123 Penghao Jiang, Ke Xin, Zifeng Wang, and Chunxi Li. Invariant meta learning for out-of-distribution
1124 generalization. *arXiv preprint arXiv:2301.11779*, 2023.

- 1134 Jean Kaddour. Stop wasting my time! saving days of imagenet and BERT training with latest
1135 weight averaging. In *Has it Trained Yet? NeurIPS 2022 Workshop*, 2022. URL [https://](https://openreview.net/forum?id=00rABUHZuz)
1136 openreview.net/forum?id=00rABUHZuz.
1137
- 1138 Jean Kaddour, Steindor Saemundsson, and Marc Deisenroth (he/him). Probabilistic Active Meta-
1139 Learning. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin (eds.), *Ad-*
1140 *vances in Neural Information Processing Systems*, volume 33, pp. 20813–20822. Curran As-
1141 sociates, Inc., 2020. URL [https://proceedings.neurips.cc/paper/2020/file/](https://proceedings.neurips.cc/paper/2020/file/ef0d17b3bdb4ee2aa741ba28c7255c53-Paper.pdf)
1142 [ef0d17b3bdb4ee2aa741ba28c7255c53-Paper.pdf](https://proceedings.neurips.cc/paper/2020/file/ef0d17b3bdb4ee2aa741ba28c7255c53-Paper.pdf).
- 1143 Jean Kaddour, Yuchen Zhu, Qi Liu, Matt J Kusner, and Ricardo Silva. Causal effect inference for
1144 structured treatments. *Advances in Neural Information Processing Systems*, 34:24841–24854,
1145 2021.
- 1146 Jean Kaddour, Linqing Liu, Ricardo Silva, and Matt Kusner. When do flat minima optimizers work?
1147 In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in*
1148 *Neural Information Processing Systems*, 2022a. URL [https://openreview.net/forum?](https://openreview.net/forum?id=vDeh2yxTvuh)
1149 [id=vDeh2yxTvuh](https://openreview.net/forum?id=vDeh2yxTvuh).
1150
- 1151 Jean Kaddour, Aengus Lynch, Qi Liu, Matt J. Kusner, and Ricardo Silva. Causal machine learning:
1152 A survey and open problems. *arXiv preprint arXiv:2206.15475*, 2022b. URL [https://arxiv.](https://arxiv.org/abs/2206.15475)
1153 [org/abs/2206.15475](https://arxiv.org/abs/2206.15475).
- 1154 Priyatham Kattakinda and Soheil Feizi. Focus: Familiar objects in common and uncommon settings,
1155 2022.
1156
- 1157 Polina Kirichenko, Pavel Izmailov, and Andrew Gordon Wilson. Last layer re-training is sufficient
1158 for robustness to spurious correlations. *arXiv preprint arXiv:2204.02937*, 2022.
1159
- 1160 Polina Kirichenko, Pavel Izmailov, and Andrew Gordon Wilson. Last layer re-training is sufficient
1161 for robustness to spurious correlations, 2023.
- 1162 Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Bal-
1163 subramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, et al. Wilds: A
1164 benchmark of in-the-wild distribution shifts. In *International Conference on Machine Learning*,
1165 pp. 5637–5664. PMLR, 2021.
- 1166 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep
1167 convolutional neural networks. In F. Pereira, C. J. Burges, L. Bottou, and K. Q. Wein-
1168 berger (eds.), *Advances in Neural Information Processing Systems*, volume 25. Curran Asso-
1169 ciates, Inc., 2012. URL [https://proceedings.neurips.cc/paper_files/paper/](https://proceedings.neurips.cc/paper_files/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf)
1170 [2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf).
1171
- 1172 David Krueger, Ethan Caballero, Joern-Henrik Jacobsen, Amy Zhang, Jonathan Binas, Dinghuai
1173 Zhang, Remi Le Priol, and Aaron Courville. Out-of-Distribution Generalization via Risk Extrapo-
1174 lation (REx). *arXiv e-prints*, art. arXiv:2003.00688, March 2020. doi: 10.48550/arXiv.2003.00688.
- 1175 David Krueger, Ethan Caballero, Joern-Henrik Jacobsen, Amy Zhang, Jonathan Binas, Dinghuai
1176 Zhang, Remi Le Priol, and Aaron Courville. Out-of-distribution generalization via risk extrapo-
1177 lation (rex). In *International Conference on Machine Learning*, pp. 5815–5826. PMLR, 2021.
1178
- 1179 Sören R Künzel, Jasjeet S Sekhon, Peter J Bickel, and Bin Yu. Metalearners for estimating heteroge-
1180 neous treatment effects using machine learning. *Proceedings of the national academy of sciences*,
1181 116(10):4156–4165, 2019.
- 1182 Tyler LaBonte, Vidya Muthukumar, and Abhishek Kumar. Dropout disagreement: A recipe for group
1183 robustness with fewer annotations. In *NeurIPS 2022 Workshop on Distribution Shifts: Connecting*
1184 *Methods and Applications*, 2022.
1185
- 1186 Yoonho Lee, Annie S. Chen, Fahim Tajwar, Ananya Kumar, Huaxiu Yao, Percy Liang, and Chelsea
1187 Finn. Surgical Fine-Tuning Improves Adaptation to Distribution Shifts, March 2023. URL
<http://arxiv.org/abs/2210.11466>. arXiv:2210.11466 [cs].

- 1188 Da Li, Jianshu Zhang, Yongxin Yang, Cong Liu, Yi-Zhe Song, and Timothy M Hospedales. Episodic
1189 training for domain generalization. In *Proceedings of the IEEE/CVF International Conference on*
1190 *Computer Vision*, pp. 1446–1455, 2019.
- 1191 Haoliang Li, Sinno Jialin Pan, Shiqi Wang, and Alex C Kot. Domain generalization with adversarial
1192 feature learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
1193 pp. 5400–5409, 2018.
- 1194 Zhiheng Li, Anthony Hoogs, and Chenliang Xu. Discover and mitigate unknown biases with
1195 debiasing alternate networks, 2022.
- 1196 Zhiheng Li, Ivan Evtimov, Albert Gordo, Caner Hazirbas, Tal Hassner, Cristian Canton Ferrer,
1197 Chenliang Xu, and Mark Ibrahim. A whac-a-mole dilemma: Shortcuts come in multiples where
1198 mitigating one amplifies others, 2023.
- 1199 Weixin Liang and James Zou. Metashift: A dataset of datasets for evaluating contextual distribution
1200 shifts and training conflicts. In *International Conference on Learning Representations*, 2022. URL
1201 <https://openreview.net/forum?id=MTex8qKavoS>.
- 1202 Evan Z Liu, Behzad Haghgoo, Annie S Chen, Aditi Raghunathan, Pang Wei Koh, Shiori Sagawa,
1203 Percy Liang, and Chelsea Finn. Just train twice: Improving group robustness without training
1204 group information. In *International Conference on Machine Learning*, pp. 6781–6792. PMLR,
1205 2021.
- 1206 Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In
1207 *Proceedings of the IEEE international conference on computer vision*, pp. 3730–3738, 2015.
- 1208 Mingsheng Long, Han Zhu, Jianmin Wang, and Michael I Jordan. Unsupervised domain adaptation
1209 with residual transfer networks. *Advances in neural information processing systems*, 29, 2016.
- 1210 Aengus Lynch, Jean Kaddour, and Ricardo Silva. Evaluating the impact of geometric and
1211 statistical skews on out-of-distribution generalization performance. In *NeurIPS 2022 Work-*
1212 *shop on Distribution Shifts: Connecting Methods and Applications*, 2022. URL <https://openreview.net/forum?id=wpT79coXAu>.
- 1213 Divyat Mahajan, Shruti Tople, and Amit Sharma. Domain generalization using causal matching. In
1214 *International Conference on Machine Learning*, pp. 7313–7324. PMLR, 2021.
- 1215 Raghav Mehta, Vítor Albiero, Li Chen, Ivan Evtimov, Tamar Glaser, Zhiheng Li, and Tal Hass-
1216 ner. You only need a good embeddings extractor to fix spurious correlations. *arXiv preprint*
1217 *arXiv:2212.06254*, 2022.
- 1218 Mazda Moayeri, Phillip Pope, Yogesh Balaji, and Soheil Feizi. A comprehensive study of image
1219 classification model sensitivity to foregrounds, backgrounds, and visual attributes, 2022a.
- 1220 Mazda Moayeri, Sahil Singla, and Soheil Feizi. Hard imagenet: Segmentations for objects with
1221 strong spurious cues. *Advances in Neural Information Processing Systems*, 35:10068–10077,
1222 2022b.
- 1223 Krikamol Muandet, David Balduzzi, and Bernhard Schölkopf. Domain generalization via invariant
1224 feature representation. In *International conference on machine learning*, pp. 10–18. PMLR, 2013.
- 1225 Kevin P. Murphy. *Probabilistic Machine Learning: An introduction*. MIT Press, 2022. URL
1226 probml.ai.
- 1227 Vaishnavh Nagarajan, Anders Andreassen, and Behnam Neyshabur. Understanding the failure modes
1228 of out-of-distribution generalization, 2020. URL <https://arxiv.org/abs/2010.15775>.
- 1229 Yannic Neuhaus, Maximilian Augustin, Valentyn Boreiko, and Matthias Hein. Spurious features
1230 everywhere – large-scale detection of harmful spurious features in imagenet, 2022. URL <https://arxiv.org/abs/2212.04871>.
- 1231 Xinkun Nie and Stefan Wager. Quasi-oracle estimation of heterogeneous treatment effects. *Biometrika*,
1232 108(2):299–319, 2021.

- 1242 NLP Connect. vit-gpt2-image-captioning (revision 0e334c7), 2022. URL [https://](https://huggingface.co/nlpconnect/vit-gpt2-image-captioning)
1243 huggingface.co/nlpconnect/vit-gpt2-image-captioning.
1244
- 1245 Daniel Omeiza, Skyler Speakman, Celia Cintas, and Komminist Weldermariam. Smooth grad-cam++:
1246 An enhanced inference level visualization technique for deep convolutional neural network models,
1247 2019.
- 1248 Karl Pearson. Mathematical contributions to the theory of evolution.—on a form of spurious
1249 correlation which may arise when indices are used in the measurement of organs. *Proceedings of*
1250 *the royal society of london*, 60(359-367):489–498, 1897.
- 1251 Jonas Peters, Peter Bühlmann, and Nicolai Meinshausen. Causal inference by using invariant
1252 prediction: identification and confidence intervals. *Journal of the Royal Statistical Society Series*
1253 *B: Statistical Methodology*, 78(5):947–1012, 2016.
- 1254
- 1255 Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of causal inference: foundations*
1256 *and learning algorithms*. The MIT Press, 2017.
- 1257
- 1258 Alexandre Rame, Corentin Dancette, and Matthieu Cord. Fishr: Invariant gradient variances for out-
1259 of-distribution generalization. In *International Conference on Machine Learning*, pp. 18347–18377.
1260 PMLR, 2022a.
- 1261 Alexandre Rame, Matthieu Kirchmeyer, Thibaud Rahier, Alain Rakotomamonjy, Patrick Gallinari,
1262 and Matthieu Cord. Diverse weight averaging for out-of-distribution generalization. *arXiv preprint*
1263 *arXiv:2205.09739*, 2022b.
- 1264
- 1265 Robin Rombach and Patrick Esser. License - a Hugging Face Space by CompVis, 2022. URL
1266 <https://huggingface.co/spaces/CompVis/stable-diffusion-license>.
- 1267
- 1268 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
1269 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Confer-*
1270 *ence on Computer Vision and Pattern Recognition*, pp. 10684–10695, 2022.
- 1271
- 1272 Elan Rosenfeld, Pradeep Ravikumar, and Andrej Risteski. Domain-adjusted regression or: Erm
1273 may already learn features sufficient for out-of-distribution generalization. *arXiv preprint*
1274 *arXiv:2202.06856*, 2022.
- 1275
- 1276 Shiori Sagawa, Pang Wei Koh, Tatsunori B. Hashimoto, and Percy Liang. Distributionally robust
1277 neural networks for group shifts: On the importance of regularization for worst-case generalization,
1278 2019a. URL <https://arxiv.org/abs/1911.08731>.
- 1279
- 1280 Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust
1281 neural networks for group shifts: On the importance of regularization for worst-case generalization.
1282 *arXiv preprint arXiv:1911.08731*, 2019b.
- 1283
- 1284 Shiori Sagawa, Aditi Raghunathan, Pang Wei Koh, and Percy Liang. An investigation of why
1285 overparameterization exacerbates spurious correlations. In *International Conference on Machine*
1286 *Learning*, pp. 8346–8356. PMLR, 2020.
- 1287
- 1288 Pamela Samuelson. Generative AI meets copyright. *Science*, 381(6654):158–161, July 2023. doi:
1289 10.1126/science.adi0656. URL [https://www.science.org/doi/10.1126/science.](https://www.science.org/doi/10.1126/science.adi0656)
1290 [adi0656](https://www.science.org/doi/10.1126/science.adi0656). Publisher: American Association for the Advancement of Science.
- 1291
- 1292 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi
1293 Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski,
1294 Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev.
1295 Laion-5b: An open large-scale dataset for training next generation image-text models, 2022. URL
<https://arxiv.org/abs/2210.08402>.

- 1296 Yuge Shi, Jeffrey Seely, Philip H. S. Torr, N. Siddharth, Awni Hannun, Nicolas Usunier, and Gabriel
1297 Synnaeve. Gradient matching for domain generalization. *arXiv preprint arXiv:2104.09937*, 2021.
1298
- 1299 Herbert A Simon. Spurious correlation: A causal interpretation. *Journal of the American statistical*
1300 *Association*, 49(267):467–479, 1954.
- 1301 Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks:
1302 Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*, 2013.
1303
- 1304 Sahil Singla and Soheil Feizi. Salient imagenet: How to discover spurious features in deep learning?,
1305 2022.
- 1306 Nimit S. Sohoni, Jared A. Dunnmon, Geoffrey Angus, Albert Gu, and Christopher Ré. No subclass
1307 left behind: Fine-grained robustness in coarse-grained classification problems, 2022.
- 1309 Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Under-
1310 standing and Mitigating Copying in Diffusion Models, May 2023. URL [http://arxiv.org/](http://arxiv.org/abs/2305.20086)
1311 [abs/2305.20086](http://arxiv.org/abs/2305.20086). arXiv:2305.20086 [cs].
1312
- 1313 Baochen Sun and Kate Saenko. Deep coral: Correlation alignment for deep domain adaptation. In
1314 *Computer Vision–ECCV 2016 Workshops: Amsterdam, The Netherlands, October 8–10 and 15–16,*
1315 *2016, Proceedings, Part III 14*, pp. 443–450. Springer, 2016.
- 1316 Damien Teney, Ehsan Abbasnejad, Simon Lucey, and Anton Van den Hengel. Evading the simplicity
1317 bias: Training a diverse set of models discovers solutions with superior ood generalization. In
1318 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.
1319 16761–16772, 2022.
- 1320
- 1321 Christopher T. H Teo and Ngai-Man Cheung. Measuring fairness in generative models, 2021.
- 1322
- 1323 Vladimir Vapnik. Principles of risk minimization for learning theory. *Advances in neural information*
1324 *processing systems*, 4, 1991.
- 1325 Joshua Vendrow, Saachi Jain, Logan Engstrom, and Aleksander Madry. Dataset interfaces: Diagnos-
1326 ing model failures using controllable counterfactual generation. *arXiv preprint arXiv:2302.07865*,
1327 2023.
- 1328
- 1329 Zhenyi Wang, Tiehang Duan, Le Fang, Qiuling Suo, and Mingchen Gao. Meta learning on a sequence
1330 of imbalanced domains with difficulty awareness. In *Proceedings of the IEEE/CVF International*
1331 *Conference on Computer Vision*, pp. 8947–8957, 2021.
- 1332
- 1333 Olivia Wiles, Sven Gowal, Florian Stimberg, Sylvestre Alvisé-Rebuffi, Ira Ktena, Krishnamurthy
1334 Dvijotham, and Taylan Cemgil. A fine-grained analysis on distribution shift. *arXiv preprint*
1335 *arXiv:2110.11328*, 2021.
- 1336
- 1337 Olivia Wiles, Isabela Albuquerque, and Sven Gowal. Discovering bugs in vision models using
1338 off-the-shelf image generation and captioning. *arXiv preprint arXiv:2208.08831*, 2022.
- 1339
- 1340 Mitchell Wortsman, Gabriel Ilharco, Samir Ya Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes,
1341 Ari S Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, et al. Model
1342 soups: averaging weights of multiple fine-tuned models improves accuracy without increasing
1343 inference time. In *International Conference on Machine Learning*, pp. 23965–23998. PMLR,
1344 2022.
- 1345
- 1346 Kai Xiao, Logan Engstrom, Andrew Ilyas, and Aleksander Madry. Noise or signal: The role of image
1347 backgrounds in object recognition. *arXiv preprint arXiv:2006.09994*, 2020.
- 1348
- 1349 Minghao Xu, Jian Zhang, Bingbing Ni, Teng Li, Chengjie Wang, Qi Tian, and Wenjun Zhang.
Adversarial domain adaptation with domain mixup, 2019.
- 1348
- 1349 Tongkun Xu, Weihua Chen, Pichao Wang, Fan Wang, Hao Li, and Rong Jin. Cdtrans: Cross-domain
transformer for unsupervised domain adaptation. *arXiv preprint arXiv:2109.06165*, 2021.

- 1350 Huaxiu Yao, Yu Wang, Sai Li, Linjun Zhang, Weixin Liang, James Zou, and Chelsea Finn. Im-
1351 proving Out-of-Distribution Robustness via Selective Augmentation. In *Proceedings of the 39th*
1352 *International Conference on Machine Learning*, pp. 25407–25437. PMLR, June 2022. URL
1353 <https://proceedings.mlr.press/v162/yao22b.html>. ISSN: 2640-3498.
- 1354 Nanyang Ye, Kaican Li, Haoyue Bai, Runpeng Yu, Lanqing Hong, Fengwei Zhou, Zhenguo Li,
1355 and Jun Zhu. Ood-bench: Quantifying and understanding two dimensions of out-of-distribution
1356 generalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
1357 *Recognition*, pp. 7947–7958, 2022.
- 1358 Yuwei Yin, Jean Kaddour, Xiang Zhang, Yixin Nie, Zhenguang Liu, Lingpeng Kong, and Qi Liu.
1359 Ttida: Controllable generative data augmentation via text-to-text and text-to-image models, 2023.
- 1360 Marvin Zhang, Henrik Marklund, Abhishek Gupta, Sergey Levine, and Chelsea Finn. Adaptive risk
1361 minimization: A meta-learning approach for tackling group shift. *arXiv preprint arXiv:2007.02931*,
1362 8:9, 2020.
- 1363 Xingxuan Zhang, Yue He, Tan Wang, Jiabin Qi, Han Yu, Zimu Wang, Jie Peng, Renzhe Xu, Zheyang
1364 Shen, Yulei Niu, et al. Nico challenge: Out-of-distribution generalization for image recognition
1365 challenges. In *European Conference on Computer Vision*, pp. 433–450. Springer, 2023.
- 1366 Shengjia Zhao, Hongyu Ren, Arianna Yuan, Jiaming Song, Noah Goodman, and Stefano Ermon.
1367 Bias and generalization in deep generative models: An empirical study, 2018.
- 1368 Evan Zheran Liu, Behzad Haghgoo, Annie S. Chen, Aditi Raghunathan, Pang Wei Koh, Shiori
1369 Sagawa, Percy Liang, and Chelsea Finn. Just Train Twice: Improving Group Robustness without
1370 Training Group Information. *arXiv e-prints*, art. arXiv:2107.09044, July 2021. doi: 10.48550/
1371 arXiv.2107.09044.
- 1372 Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep
1373 features for discriminative localization, 2015.

1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

A APPENDIX

You may include other additional sections here.