

# SUPERVISED CONTRASTIVE BLOCK DISENTANGLEMENT

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Real-world datasets often combine data collected under different experimental conditions. Although this yields larger datasets, it also introduces spurious correlations that make it difficult to accurately model the phenomena of interest. We address this by learning two blocks of latent variables to independently represent the phenomena of interest and the spurious correlations. The former are correlated with the target variable  $y$  and invariant to the environment variable  $e$ , while the latter depend on  $e$ . The invariance of the phenomena of interest to  $e$  is highly sought-after but difficult to achieve on real-world datasets. Our primary contribution is an algorithm called Supervised Contrastive Block Disentanglement (SCBD) that is highly effective at enforcing this invariance. It is based purely on supervised contrastive learning, and scales to real-world data better than existing approaches. We empirically validate SCBD on two challenging problems. The first is domain generalization, where we achieve strong performance on a synthetic dataset, as well as on Camelyon17-WILDS. SCBD introduces a single hyperparameter  $\alpha$  that controls the degree of invariance to  $e$ . When we increase  $\alpha$  to strengthen the degree of invariance, there is a monotonic improvement in out-of-distribution performance at the expense of in-distribution performance. The second is a scientific problem of batch correction. Here, we demonstrate the utility of SCBD by learning representations of single-cell perturbations from 26 million Optical Pooled Screening images that are nearly free of technical artifacts induced by the variation across wells.

## 1 INTRODUCTION

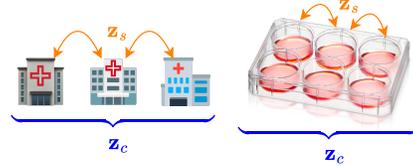
Real-world machine learning (ML) datasets often combine data collected under different experimental conditions, such as medical images or stained histopathology sections collected at different hospitals (Bánda et al., 2019; McKinney et al., 2020). This practice yields larger datasets, but the different experimental conditions alter the images’ appearance, and induce spurious correlations that make it difficult to model the phenomena of interest. While human perception is relatively robust (Makino et al., 2022b), ML models tend to rely on hospital-specific spurious correlations, and fail to generalize out-of-distribution to unseen hospitals (Koh et al., 2021).

Similar spurious correlations are a long-standing problem in experimental biology (Chandrasekaran et al., 2024), where they are called batch effects (Leek et al., 2010). They can arise between experiments conducted in different labs, within the same lab, and even within a single large parallelized experiment. Removing batch effects by batch correction is a highly-active research direction (Arevalo et al., 2024). Batch effects have been ubiquitously observed across various high-throughput lab techniques. For example, consider image-based perturbation screening, where the goal is to understand the effect of perturbations by comparing the images of perturbed cells to those from a control condition. When the images also vary due to differences in experimental conditions, this confounds our understanding of the effect of the perturbations.

In some cases, we have prior knowledge of the spurious correlations, and can take steps to remove them. For example, color-based data augmentation can remove the staining variation in histopathology images, yielding significant improvements in out-of-distribution generalization (Nguyen et al., 2023). Similarly, in experimental biology, there are manually-engineered post-processing methods that remove specific known batch effects (Carpenter et al., 2006). However, such manual approaches

054 have two important limitations. First, they require manual post-hoc quality checks to ensure the  
 055 post-processing did not remove variations of interest. Second, some spurious correlations may be  
 056 unknown, and therefore uncorrected. This motivates the development of automated approaches that  
 057 maximize the removal of the spurious correlations, while minimizing the impact on the phenomena  
 058 of interest.

059 To address these issues, we propose to learn representations of the data using two blocks of latent variables, one  
 060 encoding the variation of interest, and the other encoding the spurious correlations. We break symmetry between  
 061 the blocks by exploiting the supervisory signal of the target variable  $y$  and the environment variable  $e$ . Let  $\mathbf{x}$  be  
 062 the observation, such as a histopathology image. Let  $y$  represent the phenomenon of interest, such as the presence  
 063 of disease. Finally, let  $e$  represent the experimental conditions, such as the hospital that processed the image.  
 064 From these observed variables, we learn two embeddings  $\mathbf{z}_c$  and  $\mathbf{z}_s$ , where  $\mathbf{z}_c$  represents the variation of  $\mathbf{x}$  induced  
 065 by  $y$ , and  $\mathbf{z}_s$  represents the variation of  $\mathbf{x}$  induced by  $e$ . As we discuss in Section 2, we can also let  $\mathbf{z}_s$  represent  
 066 the variation of  $\mathbf{x}$  induced by both  $y$  and  $e$ . Our goal is to *block disentangle*  $\mathbf{z}_c$  and  $\mathbf{z}_s$  so that they independently  
 067 represent distinct information.



068 Figure 1: Spurious correlations emerge when collecting medical images from different hospitals, or conducting  
 069 single-cell perturbation screens across multiple wells.  $\mathbf{z}_s$  models these spurious correlations, while  $\mathbf{z}_c$  models the  
 070 environment-invariant correlations.

071 The promise of estimating  $\mathbf{z}_c$  such as it represents the variation of  $\mathbf{x}$  due to  $y$  in way that is invariant  
 072 to  $e$  is significant for many high-impact downstream applications. However, existing methods  
 073 for this task require additional regularization and hyperparameter tuning to ensure that  $\mathbf{z}_c$  remains  
 074 invariant with respect to  $e$ . Optimizing those hyperparameters in the presence of distribution shifts  
 075 has proven challenging in practice (Gulrajani & Lopez-Paz, 2021). While a few existing approaches  
 076 have shown success in simplified settings (Peters et al., 2016; Ganin et al., 2016; Louizos et al., 2016;  
 077 Lopez et al., 2018; Arjovsky et al., 2019; Lu et al., 2021; Kong et al., 2022), most methods tested  
 078 on real-world data have not outperformed carefully-tuned supervised learning baselines (Gulrajani  
 079 & Lopez-Paz, 2021). Consequently, the problem of learning block-disentangled representations re-  
 080 mains largely unsolved.

081 Our primary contribution is an algorithm called Supervised Contrastive Block Disentanglement  
 082 (SCBD). We claim that SCBD achieves the desired invariance to  $e$  with minimal and interpretable  
 083 hyperparameter tuning. Unlike prior work on block disentanglement that use variational or adver-  
 084 sarial objectives, our algorithm is based purely on Supervised Contrastive Learning (SCL) (Khosla  
 085 et al., 2020). Following the authors’ notation, we learn two encoder networks  $\text{Enc}_c(\cdot)$  and  $\text{Enc}_s(\cdot)$   
 086 that map  $\mathbf{x}$  to the representations given by

$$091 \quad \mathbf{r}_c := \text{Enc}_c(\mathbf{x}) \in \mathbb{R}^{D_{r_c}}, \quad \mathbf{r}_s := \text{Enc}_s(\mathbf{x}) \in \mathbb{R}^{D_{r_s}}.$$

092 We additionally learn two projection networks  $\text{Proj}_c(\cdot)$  and  $\text{Proj}_s(\cdot)$  that map the representations  
 093 to the lower-dimensional embeddings given by

$$094 \quad \mathbf{z}_c := \text{Proj}_c(\mathbf{r}_c) \in \mathbb{R}^{D_{z_c}}, \quad \mathbf{z}_s := \text{Proj}_s(\mathbf{r}_s) \in \mathbb{R}^{D_{z_s}},$$

095 which are normalized to the unit hypersphere. Finally, we learn a decoder  $\text{Dec}(\mathbf{z}_c, \mathbf{z}_s)$  which recon-  
 096 structs  $\mathbf{x}$  from  $\mathbf{z}_c$  and  $\mathbf{z}_s$ . The optimization objective consists of four losses, and is given by

$$097 \quad \min \mathcal{L}_{\mathbf{z}_c, y}^{\text{sup}} + \mathcal{L}_{\mathbf{z}_s, e}^{\text{sup}} + \alpha \mathcal{L}_{\mathbf{z}_c, e}^{\text{inv}} - \log p(\mathbf{x} | \text{Dec}(\mathbf{z}_c, \mathbf{z}_s)). \quad (1)$$

098 The first loss directly applies SCL to cluster  $\mathbf{z}_c$  with respect to  $y$ . Similarly, the second loss directly  
 099 applies SCL to cluster  $\mathbf{z}_s$  with respect to  $e$ . The third loss is our novel invariance loss, which is also  
 100 based on SCL, and ensures that  $\mathbf{z}_c$  is well-mixed with respect to  $e$ . In other words, our invariance  
 101 loss purges  $\mathbf{z}_c$  of the influence of  $e$ . The fourth loss is a reconstruction loss. We describe these losses  
 102 in detail in Section 2. SCBD incorporates a single hyperparameter  $\alpha \in \mathbb{R}_{\geq 0}$  to adjust the degree to  
 103 which  $\mathbf{z}_c$  is invariant to  $e$ . When we increase  $\alpha$ , we observe a monotonic improvement on several  
 104 downstream evaluation metrics that benefit from block disentanglement.

105 We empirically validate SCBD on three datasets spanning two challenging real-world problems.  
 106 The first problem is domain generalization (Blanchard et al., 2011; Muandet et al., 2013), where  
 107

$\mathbf{z}_c$  represents features whose correlation with  $y$  is invariant to  $e$ . We use SCBD to generalize out-of-distribution on the synthetic Colored MNIST (CMNIST) dataset, as well as the real-world histopathology dataset Camelyon17-WILDS (Koh et al., 2021). We demonstrate that SCBD enables precise control over the trade-off between in-distribution and out-of-distribution generalization performance through adjustment of the hyperparameter  $\alpha$ . Additionally, we show that on both datasets, SCBD achieves better out-of-distribution performance relative to the conventional baselines in the literature.

The second problem is batch correction, where we apply SCBD to a dataset of images of over 26 million individual cells (Funk et al., 2022). The cells are treated with 5,050 genetic perturbations which are labeled as  $y$ , and collected across 34 wells which are labeled as  $e$ . We use SCBD to represent the effect of the perturbation with  $\mathbf{z}_c$ , and the variation across wells with  $\mathbf{z}_s$ . We show that relative to strong baselines, SCBD provides estimates of  $\mathbf{z}_c$  that preserves biological signal and are significantly less sensitive to batch effects.

## 2 SUPERVISED CONTRASTIVE BLOCK DISENTANGLEMENT

We now define the individual terms in the SCBD optimization objective in Equation 1. Our starting point is a probabilistic interpretation of SCL that helps derive our novel invariance loss. Following the notation from Khosla et al. (2020), let  $I$  be the set of indices of examples within a minibatch. For each anchor point  $i \in I$ , we denote the set of the remaining examples as  $\mathcal{A}(i) = I \setminus \{i\}$ . In SCL, anchor points are compared to other examples via their dot product. We define  $|\mathcal{A}(i)|$  independent Bernoulli random variables  $M_{i,c}^j$  for  $j$  in  $\{1, \dots, |\mathcal{A}(i)|\}$  to represent whether  $\mathbf{z}_c^i$  is matched with  $\mathbf{z}_c^j$ . The matching probability is defined as

$$P(M_{i,c}^j = 1) = \frac{\exp(\mathbf{z}_c^i \cdot \mathbf{z}_c^j / \tau)}{\sum_{a \in \mathcal{A}(i)} \exp(\mathbf{z}_c^i \cdot \mathbf{z}_c^a / \tau)}.$$

The notion of matching is relative to the rest of the examples in  $\mathcal{A}(i)$ , which is why we use a softmax despite the random variable being binary. A similar definition holds for the random variable  $M_{i,s}^j$ , defined in the space of  $\mathbf{z}_s$ .

The first term in Equation 1 is a direct application of SCL, and is given by

$$\mathcal{L}_{\mathbf{z}_c, y}^{\text{sup}} = - \sum_{i \in I} \frac{1}{|\mathcal{P}_y(i)|} \sum_{p \in \mathcal{P}_y(i)} \log P(M_{i,c}^p = 1),$$

where  $\mathcal{P}_y(i) = \{j \in \mathcal{A}(i) : y^i = y^j\}$  are the positive pairs for the anchor point  $i$  with respect to  $y$ . This represents the negative log joint probability of observing the positive pairs, normalized by the number of positive pairs, and summed across all anchor points. Minimizing this loss clusters  $\mathbf{z}_c$  with respect to  $y$ .

The second term in Equation 1 is also a direct application of SCL, and is given by

$$\mathcal{L}_{\mathbf{z}_s, e}^{\text{sup}} = - \sum_{i \in I} \frac{1}{|\mathcal{P}_e(i)|} \sum_{p \in \mathcal{P}_e(i)} \log P(M_{i,s}^p = 1),$$

where  $\mathcal{P}_e(i) = \{j \in \mathcal{A}(i) : e^i = e^j\}$  are the positive pairs for the anchor point  $i$  with respect to  $e$ . Minimizing this loss clusters  $\mathbf{z}_s$  with respect to  $e$ . As we later discuss in our experiments (Section 4.1), it can be useful to let  $\mathbf{z}_s$  represent the variation of  $\mathbf{x}$  with respect to the pair  $(y, e)$ , rather than just  $e$ . We can do this by replacing  $\mathcal{L}_{\mathbf{z}_s, e}^{\text{sup}}$  with

$$\mathcal{L}_{\mathbf{z}_s, (y, e)}^{\text{sup}} = - \sum_{i \in I} \frac{1}{|\mathcal{P}_{(y, e)}(i)|} \sum_{p \in \mathcal{P}_{(y, e)}(i)} \log P(M_{i,s}^p = 1),$$

where  $\mathcal{P}_{(y, e)}(i) = \{j \in \mathcal{A}(i) : y^i = y^j, e^i = e^j\}$  are the positive pairs for the anchor point  $i$  with respect to the pairs of labels  $(y, e)$ .

The third term in Equation 1 is our novel invariance loss. We define  $\mathcal{N}_e(i) = \mathcal{A}(i) \setminus \mathcal{P}_e(i)$  as the negative pairs with respect to the label  $e$ . Since  $\{\mathcal{P}_e(i), \mathcal{N}_e(i)\}$  is a partition of  $\mathcal{A}(i)$ , we consider

the binary classification task of whether  $\mathbf{z}_c^i$  is more likely to be matched with its positive or negative pairs with respect to  $e$ . One way to perform this classification is to predict that  $i$  is more likely to be matched with  $\mathcal{P}_e(i)$  if

$$\sum_{p \in \mathcal{P}_e(i)} \log P(M_{i,c}^p = 1) > \sum_{n \in \mathcal{N}_e(i)} \log P(M_{i,c}^n = 1).$$

Since our goal is to make  $\mathbf{z}_c$  invariant to  $e$ , we optimize  $\mathbf{z}_c$  to make this classifier fail. We do this by minimizing

$$\mathcal{L}_{\mathbf{z}_c, e}^{\text{inv}} = \left| \sum_{p \in \mathcal{P}_e(i)} \log P(M_{i,c}^p = 1) - \sum_{n \in \mathcal{N}_e(i)} \log P(M_{i,c}^n = 1) \right|,$$

which makes it equally probable that  $\mathbf{z}_c^i$  is matched with its positive and negative pairs with respect to  $e$ . In other words, it makes  $\mathbf{z}_c^i$  well-mixed with respect to  $e$ . This is analogous to adversarial approaches that train a discriminator to predict  $e$ , where the goal of representation learning is to fool the discriminator (Ganin et al., 2016; Edwards & Storkey, 2016). However, it can be difficult to apply these adversarial methods due to the complexity of minimax optimization. SCBD circumvents the need to train a discriminator, since the dot products between pairs of  $\mathbf{z}_c$  can be used to predict  $e$ .

The fourth term in Equation 1 reconstructs  $\mathbf{x}$  from  $\mathbf{z}_c$  and  $\mathbf{z}_s$ . Importantly, we only use the reconstruction loss to optimize the decoder parameters, while holding  $\mathbf{z}_c$  and  $\mathbf{z}_s$  fixed. Therefore, the representation learning of  $\mathbf{z}_c$  and  $\mathbf{z}_s$  is done purely through supervised contrastive learning. It is possible to train the encoders and decoder jointly, which would likely improve the reconstruction quality. However, this adds the further complexity of balancing the relative contributions of the supervised contrastive and reconstruction losses by incorporating an additional hyperparameter. We leave this to future work, and focus on achieving strong performance on downstream tasks.

### 3 IMPROVEMENTS TO VARIATIONAL APPROACHES FALL SHORT OF SCBD

As a basis of comparison for SCBD, we develop a block disentanglement approach based on Identifiable Variational Autoencoders (iVAEs) (Khemakhem et al., 2020). Several extensions to Variational Autoencoders (VAEs) (Kingma & Welling, 2014) address the problem of invariance to auxiliary variables, including the Variational Fair Autoencoder (Louizos et al., 2016) and the HSIC-constrained VAE (Lopez et al., 2018). These methods learn a single block of latent variables, and apply additional regularization to achieve invariance to an auxiliary variable. While successful in simple settings, these approaches have not gained widespread adoption for large-scale imaging data. Wang et al. (2023) applied contrastive learning to train a VAE with two blocks of latent variables, where one block does not condition on any auxiliary variables, and the other does. Our approach differs from theirs because we do not use contrastive learning, and both of our latent blocks condition on auxiliary variables.

We specify an iVAE with the same blocks of latent variables as SCBD. The generative model is defined as

$$p_\theta(\mathbf{x}, \mathbf{z}_c, \mathbf{z}_s \mid y, e) = p_\theta(\mathbf{x} \mid \mathbf{z}_c, \mathbf{z}_s) p_\theta(\mathbf{z}_c \mid y) p_\theta(\mathbf{z}_s \mid e),$$

while the inference model is defined as

$$q_\phi(\mathbf{z}_c, \mathbf{z}_s \mid \mathbf{x}, e) = q_\phi(\mathbf{z}_c \mid \mathbf{x}) q_\phi(\mathbf{z}_s \mid \mathbf{x}, e).$$

We fit this model by maximizing the evidence lower bound (ELBO) (Jordan et al., 1999), given by

$$\begin{aligned} \min_{\theta, \phi} \mathbb{E}_{q_\phi(\mathbf{z}_c \mid \mathbf{x}) q_\phi(\mathbf{z}_s \mid \mathbf{x}, e)} [-\log p_\theta(\mathbf{x} \mid \mathbf{z}_c, \mathbf{z}_s)] \\ + D_{KL}(q_\phi(\mathbf{z}_c \mid \mathbf{x}) \parallel p_\theta(\mathbf{z}_c \mid y)) + D_{KL}(q_\phi(\mathbf{z}_s \mid \mathbf{x}, e) \parallel p_\theta(\mathbf{z}_s \mid e)). \end{aligned}$$

Empirically, conditioning the posterior of  $\mathbf{z}_s$  on both  $\mathbf{x}$  and  $e$ , rather than just  $\mathbf{x}$ , significantly impacts downstream performance. We demonstrate this using ablation studies, comparing the two versions iVAE ( $q_\phi(\mathbf{z}_s \mid \mathbf{x}, e)$ ) and iVAE ( $q_\phi(\mathbf{z}_s \mid \mathbf{x})$ ). We hypothesize that conditioning the posterior of  $\mathbf{z}_s$  on  $e$  makes it easier to encode the variation with respect to  $e$  in  $\mathbf{z}_s$ , which reduces the incentive to

216 encode it in  $\mathbf{z}_c$ . We use a mixture of experts approach to condition on  $e$ , where a neural network  
 217 takes in  $\mathbf{x}$  and outputs separate posterior parameters for each value of  $e$ .  
 218

219 Despite this improvement, we find that iVAE performs worse than SCBD in general. VAE-based  
 220 block disentanglement methods inherently struggle to balance reconstruction and KL divergence  
 221 minimization, leading to several failure modes. First, posterior collapse happens when the KL term  
 222 is trivially minimized to zero by making the latent variables uninformative (Bowman et al., 2015;  
 223 Razavi et al., 2019; Fu et al., 2019; Dai et al., 2020; Wang et al., 2021). Second, prior collapse  
 224 occurs when learned parameters for  $p_\theta(\mathbf{z}_c | y)$  collapse to the uninformative prior  $p_\theta(\mathbf{z}_c)$ . Third,  
 225 numerical instability necessitates heuristics such as gradient clipping or skipping (Child, 2021).  
 226 These issues significantly limit the ability to train VAEs with large-capacity neural networks, likely  
 227 explaining why open-source VAE implementations rarely use generic image encoders like those in  
 torchvision (Marcel & Rodriguez, 2010).  
 228

## 229 4 EXPERIMENTS

230  
 231 We empirically validate SCBD on three datasets spanning two difficult real-world problems. We  
 232 discuss domain generalization in Section 4.1, and batch correction in Section 4.2. In domain gen-  
 233 eralization, we use one synthetic and one realistic dataset, while in batch correction we use one  
 234 large-scale realistic dataset with over 26 million images. The two problems are similar in that in  
 235 both cases, we want  $\mathbf{z}_c$  to represent the correlation between  $\mathbf{x}$  and  $y$  that is invariant to  $e$ , and  $\mathbf{z}_s$   
 236 to encode the remaining spurious correlations that depend on  $e$ . The key difference between the  
 237 two problems relates to the evaluation. In domain generalization, we evaluate the ability to predict  
 238  $y$  given  $\mathbf{z}_c$  on an out-of-distribution test set. In contrast, in batch correction the evaluation is in-  
 239 distribution, and measures the degree to which  $\mathbf{z}_c$  discards the information in  $e$ , while preserving the  
 240 information in  $y$ .  
 241

### 242 4.1 DOMAIN GENERALIZATION

#### 243 4.1.1 PROBLEM DESCRIPTION

244  
 245 Domain generalization is an out-of-distribution generalization problem, where the data come from  
 246 different environments. Environments represent different conditions under which data are gener-  
 247 ated, such as the hospital that collected the samples. We assume data are sampled from a family of  
 248 distributions  $p_{\text{all}} = \{p_e(\mathbf{x}_e, y_e) : e \in \mathcal{E}_{\text{all}}\}$  indexed by the environment  $e \in \mathcal{E}_{\text{all}} \subseteq \mathbb{N}$ . The training  
 249 data are sampled from  $p_{\text{tr}} = \{p_e(\mathbf{x}_e, y_e) : e \in \mathcal{E}_{\text{tr}}\}$ , where  $\mathcal{E}_{\text{tr}} \subset \mathcal{E}_{\text{all}}$  is the set of training environ-  
 250 nments. The test data are sampled from  $p_{\text{te}} = \{p_e(\mathbf{x}_e, y_e) : e \in \mathcal{E}_{\text{te}}\}$ , where  $\mathcal{E}_{\text{te}} \subset \mathcal{E}_{\text{all}}$  is the set of test  
 251 environments. Because  $\mathcal{E}_{\text{tr}}$  and  $\mathcal{E}_{\text{te}}$  are disjoint, there is a distribution shift between  $p_{\text{tr}}$  and  $p_{\text{te}}$ . The  
 252 goal is to predict  $y$  from  $\mathbf{x}$  in a way that is invariant to  $e$ , so that we can generalize from  $p_{\text{tr}}$  to  $p_{\text{te}}$ .  
 253

#### 254 4.1.2 IN- AND OUT-OF-DISTRIBUTION PERFORMANCE MUST BE NEGATIVELY CORRELATED

255  
 256 We begin by precisely characterizing the conditions under which SCBD should be effective at do-  
 257 main generalization. This is important, as it motivates our choice of datasets that we use in our  
 258 experiments. The conditions are intuitively simple and empirically testable. SCBD requires the  
 259 existence of spurious features in the training environments where the more they are learned, in-  
 260 distribution performance improves, and out-of-distribution performance worsens. In other words,  
 261 datasets need to exhibit a trade-off between in- and out-of-distribution performance. SCBD pre-  
 262 vents the learning of such spurious features, since they are predictive of the training environments.  
 263 This promotes the learning of features that are invariant to the environment, and thus generalize on  
 264 the test environments.

265 We therefore want to evaluate SCBD on datasets that exhibit this trade-off. Fortunately, there is an  
 266 empirical test for this, which is to train ERM across a large region of the hyperparameter search  
 267 space, and check whether there are regions where in-distribution performance is strong, and is neg-  
 268 atively correlated with out-of-distribution performance. Teney et al. (2024) carried out such a study,  
 269 and found the trade-off to be particularly prominent on the Camelyon17-WILDS (Koh et al., 2021)  
 dataset. We therefore include this dataset in our experiments.

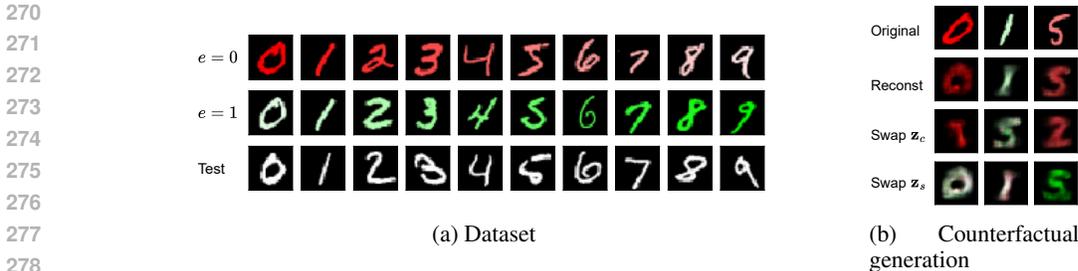


Figure 2: Colored MNIST. (a) There is an environment-dependent correlation between color and digit on the training set, which does not persist on the test set where all digits are white. (b) We can generate images counterfactually using SCBD. When we swap  $z_c$  across examples, it changes the digit without affecting the color. In contrast, when we swap  $z_s$  across examples, it changes the color without affecting the digit. By composing digit and color independently, we generate images outside of the support of the training distribution, such as a light red one (bottom middle) and a bright green five (bottom right).

This trade-off between in- and out-of-distribution performance is the exception rather than the rule for domain generalization datasets. That is, despite the datasets being constructed to have qualitatively different training and test environments, it is often the case that in- and out-of-distribution performance are positively correlated. Wenzel et al. (2022) reached this conclusion by carrying out a large-scale empirical study involving 172 datasets, including those in the DomainBed (Gulrajani & Lopez-Paz, 2021) and WILDS (Koh et al., 2021) suites. It is difficult to outperform ERM when the correlation is positive, which may explain why Gulrajani & Lopez-Paz (2021) found ERM to be state-of-the-art across the DomainBed suite.

#### 4.1.3 DATASETS

In addition to Camelyon17-WILDS, we experiment with one synthetic dataset. This dataset is called Colored MNIST (CMNIST), and extends the version from Arjovsky et al. (2019). The target label  $y \in \{0, \dots, 9\}$  represents the digit. There are two training environments and a test environment. In the training environments  $e \in \{0, 1\}$ , there is an environment-dependent correlation between the color and  $y$  (Figure 2a). For  $e = 0$  the color changes from dark to light red as the digit increases. In contrast, for  $e = 1$  the color changes from light to dark green as the digit increases. All digits are white in the test environment. This presents a severe distribution shift, since color is perfectly predictive of  $y$  in the training environments, but is unpredictable in the test environment. Details regarding the data generating process are in Appendix A.2.1. We train ERM across a range of learning rates and maximum training steps on this dataset, and observe that in- and out-of-distribution performance are negatively correlated (Appendix Figure 5). This satisfies the assumptions of SCBD, so therefore we expect it to be effective on this dataset. We expect  $z_c$  to encode the digit, and  $z_s$  to encode the environment-specific colors.

Camelyon17-WILDS (Koh et al., 2021) is a patch-based variant of the original Camelyon17 dataset (Bánci et al., 2019) of histopathology images of breast tissue, and represents a binary classification task of predicting the presence of a tumor. The data were collected in five hospitals, and have significant inter-hospital batch effects. It has been reported that for similar datasets, the most significant batch effects are from differences in how the slides are stained (Tellez et al., 2019). As mentioned previously, Teney et al. (2024) showed that this dataset exhibits a trade-off between in- and out-of-distribution performance, and therefore satisfies the assumptions of SCBD. We also verify this in Appendix Figure 11. On this dataset, we want  $z_c$  to represent the biomarkers of disease that are invariant across hospitals, and  $z_s$  to represent the hospital-specific spurious correlations.

#### 4.1.4 BASELINES

We compare SCBD to a diverse range of algorithms that are considered to be standard baselines in the domain generalization literature. This includes Empirical Risk Minimization (ERM) (Vapnik, 1995), CORrelation ALignment (CORAL) (Sun & Saenko, 2016), Domain-Adversarial Neural

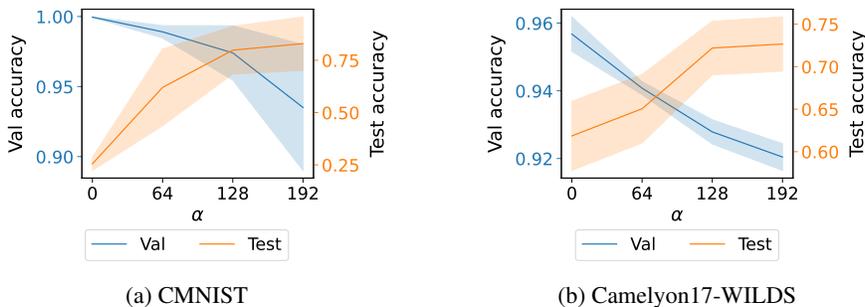


Figure 3: Increasing  $\alpha$  strengthens the degree that  $\mathbf{z}_c$  is invariant to  $e$ , and monotonically improves test accuracy at the expense of validation accuracy.

Table 1: Test accuracy (%) for domain generalization for ten random seeds. SCBD with  $\alpha = 192$  significantly outperforms all baselines on both CMNIST and Camelyon17-WILDS.

Algorithm	CMNIST	Camelyon17-WILDS
SCBD ( $\alpha = 0$ )	25.5 $\pm$ 3.0	61.9 $\pm$ 3.8
SCBD ( $\alpha = 192$ )	<b>82.9 <math>\pm</math> 12.1</b>	<b>72.7 <math>\pm</math> 3.0</b>
ERM	37.8 $\pm$ 2.6	65.8 $\pm$ 4.9
CORAL	37.6 $\pm$ 3.6	59.5 $\pm$ 7.7
DANN	39.0 $\pm$ 4.5	55.2 $\pm$ 6.7
IRM	37.0 $\pm$ 4.2	66.3 $\pm$ 2.1
Fish	48.2 $\pm$ 3.5	49.1 $\pm$ 0.9
Group DRO	35.0 $\pm$ 2.9	68.4 $\pm$ 7.3
iVAE ( $q_\phi(\mathbf{z}_s   \mathbf{x}, e)$ )	52.1 $\pm$ 37.6	52.0 $\pm$ 2.0
iVAE ( $q_\phi(\mathbf{z}_s   \mathbf{x})$ )	37.7 $\pm$ 29.4	51.9 $\pm$ 4.3

Networks (DANN) (Ganin et al., 2016), Invariant Risk Minimization (IRM) (Arjovsky et al., 2019), Fish (Shi et al., 2022), and Group Distributionally Robust Optimization (Group DRO) (Sagawa et al., 2020). We additionally compare against our two versions of iVAE from Section 3 for completeness.

#### 4.1.5 QUALITATIVE RESULTS

Our image generation results in Figure 2b qualitatively demonstrate that SCBD achieves block disentanglement. This is possible on CMNIST because we know that the ground-truth phenomenon of interest is the digit, and the spurious correlation is the color. These results show that when we swap  $\mathbf{z}_c$  between examples, it changes the digit without affecting the color. In contrast, when we swap  $\mathbf{z}_s$  between examples, it changes the color without affecting the digit. Note that the quality of the reconstructed images is relatively poor because, as mentioned in Section 2, the decoder is not trained jointly with the encoders. We leave it to future work to train the decoder jointly and improve the image reconstruction capability of SCBD. We provide similar visualization results with the iVAE in Appendix Figure 6.

#### 4.1.6 QUANTITATIVE RESULTS

We present two kinds of quantitative results. In Figure 3, we show that by increasing  $\alpha$ , SCBD removes spurious correlations that are specific to the training environments. This encourages the learning of features that are invariant to the environment, which yields a smooth and monotonic trade-off between in- and out-of-distribution performance on both datasets.

In Table 1, we show the test accuracy on both datasets for SCBD and the baseline algorithms. We report the average and standard deviation for ten random seeds. For the baseline algorithms, we optimize the hyperparameters with respect to the performance on the in-distribution validation set. The hyperparameter search space for each algorithm is provided in Appendix Table 5. Most of the baseline results for Camelyon17-WILDS are taken from the authors’ leaderboard, with the

378 exception of Fish (Shi et al., 2022), which we evaluate ourselves. Our results for Fish are weaker  
 379 than those reported on the leaderboard, because we additionally included the pretraining duration in  
 380 the hyperparameter search space. The leaderboard results used the value of this hyperparameter that  
 381 achieved the best test accuracy, as described in the appendix of Shi et al. (2022).

382 For SCBD, we apply the same model selection procedure to optimize the learning rate and weight  
 383 decay. We do not optimize  $\alpha$  during model selection, since this would result in choosing  $\alpha = 0$ .  
 384 We report the test accuracy for  $\alpha = 0$  and  $\alpha = 192$  as evidence that the invariance loss in SCBD  
 385 is effective at removing spurious correlations and improving out-of-distribution performance. With  
 386  $\alpha = 192$ , SCBD significantly outperforms all baseline algorithms across both datasets. Tuning  $\alpha$   
 387 corresponds to model selection with respect to an unknown test distribution, which is a difficult open  
 388 problem (Gulrajani & Lopez-Paz, 2021), and is a limitation shared by other works (Makino et al.,  
 389 2022a; Wortsman et al., 2022).

390 Also, we demonstrate the robustness of our approach to the choice of hyperparameters by providing  
 391 the results of ablation studies in Appendix A.2.1 and A.2.2, where we vary  $D_{z_c}$  and  $D_{z_s}$ , the batch  
 392 size, and the degree of weight decay.

393 We additionally experiment with PACS (Li et al., 2017) and VLCS (Fang et al., 2013) from Do-  
 394 mainBed (Gulrajani & Lopez-Paz, 2021), and include the results in Appendix Sections A.2.3 and  
 395 A.2.4. We find that these datasets exhibit a positive correlation between in- and out-of-distribution  
 396 performance, which is consistent with Wenzel et al. (2022). Since this violates the assumptions of  
 397 SCBD, we are unable to trade-off in- and out-of-distribution performance, as expected.

## 399 4.2 BATCH CORRECTION WITH A REAL-WORLD OPTICAL POOLED SCREEN DATASET

### 401 4.2.1 PROBLEM DESCRIPTION AND DATASET

402 Having demonstrated the efficacy of SCBD on domain generalization, we proceed to a related but  
 403 different problem called batch correction. Here, we experiment with one realistic single-cell pertur-  
 404 bation dataset that is of significantly large scale. We use the Optical Pooled Screen (OPS) (Feldman  
 405 et al., 2019) dataset from Funk et al. (2022) comprised of 26 million images of single cells, each  
 406 perturbed with one of 5,050 genetic perturbations targeting an expressed gene, including one non-  
 407 targeting control. Such data are collected in order to understand the effect of each perturbation on  
 408 cellular morphology. The  $100 \times 100$  pixel images have four channels that measure staining informa-  
 409 tion for key cellular features: DNA damage, F-actin, DNA content, and microtubules. Each channel  
 410 therefore measures a unique aspect of a cell’s phenotype, which taken together shed light on how  
 411 each perturbed genes affects the cell. An important problem in the field is to build a cartography  
 412 of perturbation effects on cells, by grouping perturbed genes by their similarity on the phenotypic  
 413 level (Celik et al., 2024). This perturbation map is then interpreted to characterize the function of  
 414 unknown genes, recapitulate protein complexes, and highlight interacting pathways (Rood et al.,  
 415 2024).

416 OPS technologies generate large quantities of data in a cost effective manner by conducting several  
 417 batches of experiments in parallel. In this case, the data were collected at a single lab but using  
 418 34 wells. There can be significant unintended variation across wells, solely based on seemingly  
 419 minor differences in experimental conditions. For example, if the wells are stained sequentially,  
 420 the difference in elapsed time can result in different image brightness across wells. Our goal with  
 421 SCBD is to capture this unintended variation across wells in  $z_s$ , so that  $z_c$  is an unconfounded  
 422 representation of the impact of genetic perturbations on cell morphology.

423 For each image of a single cell  $x$ ,  $y$  labels the genetic perturbation that was applied, and  $e$  labels  
 424 which of the 34 wells the cell was in. By optimizing the SCBD objective in Equation 1, we ensure  
 425 that the variation in the images due to the perturbation is represented by  $z_c$ , and the variation due  
 426 to the well  $e$  is represented by  $z_s$ . We can then use  $z_c$  for downstream analysis. For this task,  
 427 we are using  $z_c$  with  $D_{z_c} = D_{z_s} = 64$ , whereas we used  $r_c$  for domain generalization. This is  
 428 because all of our baselines for this task use 64 dimensional embeddings, so the lower-dimensional  
 429  $z_c$  helps ensure a fair comparison. We show results using ResNet-18 encoders in the main text, and  
 430 additionally provide results using DenseNet-121 encoders in Appendix Figure 25.

431 We evaluate two tasks to understand the degree to which we remove the influence of  $e$ , while pre-  
 serving the information in  $y$ . We describe the tasks at a high level here, and provide the details in

432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485

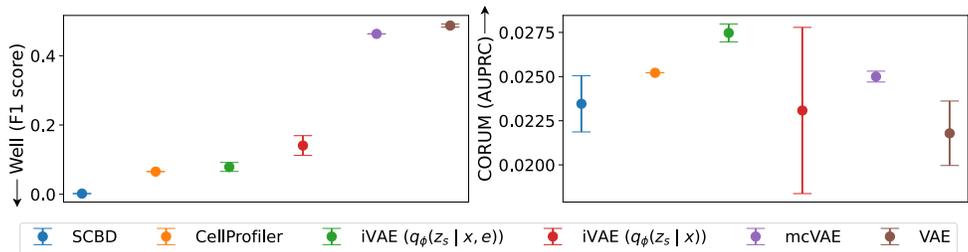


Figure 4: Comparison of SCBD to CellProfiler and VAE-based baselines on real-world batch correction. Left: Performance of predicting the well label  $e$ . Right: Performance on predicting protein complex membership (biological content). SCBD is almost completely uninformative of the well, while retaining a similar level of biological information as CellProfiler.

Appendix A.3. The first task is to use the perturbation embeddings to predict  $e$ , which measures the sensitivity to inter-well batch effects. We fit a linear classifier on top of each of the embeddings, and compute the F1 score. We want the performance on this task to be weak.

The second task is CORUM prediction, which is one measure of the biological information content in the embeddings. This task relies on the CORUM database (Ruepp et al., 2010) as the ground truth of whether two genes are functionally related based on their membership in the same protein complex (a definition previously used in the context of this biological screen). We take the biological embeddings corresponding to those genetic perturbations  $y$ , interpret their dot product as the prediction that they are similar, and use these predictions to compute the area under the precision-recall curve. Unlike the first task, we want the performance on this task to be strong.

#### 4.2.2 BASELINES

CellProfiler (Carpenter et al., 2006) is the most important baseline that we compare SCBD against. It is an open-source software that takes in an image of a cell, and outputs several thousand manually-engineered morphological features that describe the cell’s phenotype. It is a very strong baseline in which substantial human-expert effort has been invested, and its representations are post-processed to effectively remove the variation across plates and wells. Following conventional practice, we use the top-64 principal components of the full set of CellProfiler features.

The remaining baselines are all based on VAEs, which are considered conventional. We experiment with iVAE ( $q_\phi(\mathbf{z}_s | \mathbf{x}, e)$ ) and iVAE ( $q_\phi(\mathbf{z}_s | \mathbf{x})$ ) from Section 3, as well as the Multi-Contrastive VAE (mcVAE) (Wang et al., 2023), which uses two blocks of latent variables in order to represent the perturbation effect and the natural cell-to-cell variation. While it was previously shown that mcVAE is effective at modeling genetic perturbations, it has a significant weakness in that it does not effectively correct for batch effects. Finally, for our simplest baseline we use a vanilla VAE (Kingma & Welling, 2014), which has a single block of latent variables, and ignores  $y$  and  $e$ . For all VAE-based models, we use 64 dimensional latent variables in each block. The perturbation embedding is  $\mathbf{z}_c$  for SCBD and the iVAEs. For mcVAE it is the block of salient variables, and for CellProfiler and the vanilla VAE, there is only a single block of latent variables.

#### 4.2.3 QUANTITATIVE RESULTS

We show our results on both tasks in Figure 4. SCBD achieves close to zero predictive performance on the well-prediction task, while retaining a similar level of biological information as CellProfiler. Thus,  $\mathbf{z}_c$  estimated with SCBD can be used by biologists for downstream analysis, and they can be confident that any conclusions reached are not due to the inter-well variation. Our results also show that for the iVAE baselines, additionally conditioning the posterior of  $\mathbf{z}_s$  on  $e$  significantly improves the results on both tasks. Although mcVAE performs better than the vanilla VAE on CORUM due to its ability to incorporate the perturbation labels  $y$ , they are both highly susceptible to the inter-well batch effects. This highlights the fact that explicit regularization is required in order to purge the effect of  $e$  from the embeddings, and that this does not occur naturally.

## 5 RELATED WORK

### 5.1 DISENTANGLED REPRESENTATION LEARNING

Our goal of block disentanglement is closely related to that of disentangled representation learning, which assumes that a relatively small number of independent factors are sufficient to explain the important patterns of variation in  $\mathbf{x}$ . Disentangled representation learning is typically cast as learning a latent variable  $\mathbf{z} \in \mathbb{R}^{D_z}$ , where  $\mathbf{z}$  is disentangled if its individual components  $z_1, \dots, z_{D_z}$  are independent and semantically meaningful (Higgins et al., 2017; Esmaeili et al., 2019; Kim & Mnih, 2018; Chen et al., 2018). This informal definition of disentanglement is generally agreed upon, and it is not trivial to define this concept quantitatively (Eastwood & Williams, 2018; Higgins et al., 2018). This is related to independent component analysis (Comon, 1994; Jutten & Herault, 1991; Hyvärinen & Oja, 2000), which makes the additional assumption that the encoding is noiseless.

With block disentanglement, instead of assuming there are  $D_z$  independent scalar factors, we assume there are two independent vector-valued factors  $\mathbf{z}_c \in \mathbb{R}^{D_{z_c}}$  and  $\mathbf{z}_s \in \mathbb{R}^{D_{z_s}}$ . Recent works study identifiability for block disentanglement (Von Kügelgen et al., 2021; Lachapelle & Lacoste-Julien, 2022; Kong et al., 2022; Lachapelle et al., 2024; Lopez et al., 2024). While we believe this is an important research direction, we focus on developing a simple algorithm that achieves strong empirical results on difficult real-world problems.

### 5.2 INVARIANT REPRESENTATION LEARNING

The challenge of domain generalization has gained significant attention as ML systems often fail to generalize out-of-distribution. Peters et al. (2016) introduced a framework for causal inference using invariant prediction, helping maintain predictive accuracy under interventions or environmental changes. Building on this foundation, Arjovsky et al. (2019) proposed IRM, a learning paradigm for learning an embedding of the data representation such that the optimal classifier on top of that representation remains invariant across different environments. These works, as well as many extensions (Lu et al., 2021), have been benchmarked on datasets created by the research community, such as those in the DomainBed (Gulrajani & Lopez-Paz, 2021) and WILDS (Koh et al., 2021) suites. Gulrajani & Lopez-Paz (2021) revealed that with rigorous model selection, ERM often achieves state-of-the-art performance, challenging the perceived benefits of more complex domain generalization methods.

## 6 CONCLUSION

We presented Supervised Contrastive Block Disentanglement (SCBD), an algorithm for block disentanglement that is based purely on supervised contrastive learning. We use SCBD to estimate  $\mathbf{z}_c$  such that it represents the correlation between  $\mathbf{x}$  and  $y$  that is invariant to  $e$ . This invariance, which is considered difficult to achieve in practice, allows us to solve two difficult real-world problems. The first is domain generalization, where we achieve strong out-of-distribution generalization on a synthetic dataset called Colored MNIST, as well as a real-world histopathology dataset called Camelyon17-WILDS. The second is batch correction, where we use SCBD to learn representations of single-cell perturbations from over 26 million images that are nearly free of inter-well batch effects.

We believe a promising direction for future work is to investigate how to jointly train the decoder to combine the representation learning capabilities of supervised contrastive learning and generative modeling. While iVAE failed at domain generalization, it achieved strong performance on the CORUM task in the batch correction problem. We interpret this as a sign that image reconstruction can yield additional useful features that SCBD is currently not capturing. With improved generative modeling, SCBD has the potential to be used for impactful counterfactual image generation, such as generating images of the same cell under different perturbations. Also, in this work we assumed access to the variable  $e$ , which labels the source of unwanted variation. We leave it to future work to learn this variable from data.

## REFERENCES

- 540  
541  
542 John Arevalo, Ellen Su, Jessica D Ewald, Robert van Dijk, Anne E Carpenter, and Shantanu Singh.  
543 Evaluating batch correction methods for image-based cell profiling. *Nature Communications*,  
544 2024.
- 545 Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization.  
546 *arXiv*, 2019.
- 547 Péter Bándi, Oscar Geessink, Quirine Manson, Marcory Van Dijk, Maschenka Balkenhol, Meyke  
548 Hermsen, Babak Ehteshami Bejnordi, Byungjae Lee, Kyunghyun Paeng, et al. From detection  
549 of individual metastases to classification of lymph node status at the patient level: The CAME-  
550 LYON17 challenge. *IEEE Trans. Medical Imaging*, 2019.
- 552 Gilles Blanchard, Gyemin Lee, and Clayton Scott. Generalizing from several related classification  
553 tasks to a new unlabeled sample. *Advances in Neural Information Processing Systems*, 2011.
- 554 Samuel R Bowman, Luke Vilnis, Oriol Vinyals, Andrew M Dai, Rafal Jozefowicz, and Samy Ben-  
555 gio. Generating sentences from a continuous space. *arXiv*, 2015.
- 557 Anne E Carpenter, Thouis R Jones, Michael R Lamprecht, Colin Clarke, In Han Kang, Ola Friman,  
558 David A Guertin, Joo Han Chang, Robert A Lindquist, Jason Moffat, et al. Cellprofiler: image  
559 analysis software for identifying and quantifying cell phenotypes. *Genome Biology*, 2006.
- 560 Safiye Celik, Jan-Christian Hütter, Sandra Melo Carlos, Nathan H Lazar, Rahul Mohan, Conor  
561 Tillinghast, Tommaso Biancalani, Marta M Fay, Berton A Earnshaw, and Imran S Haque. Build-  
562 ing, benchmarking, and exploring perturbative maps of transcriptional and morphological data.  
563 *bioRxiv*, 2024.
- 564 Srinivas Niranj Chandrasekaran, Beth A Cimini, Amy Goodale, Lisa Miller, Maria Kost-Alimova,  
565 Nasim Jamali, John G Doench, Briana Fritchman, Adam Skepner, Michelle Melanson, et al. Three  
566 million images and morphological profiles of cells treated with matched chemical and genetic  
567 perturbations. *Nature Methods*, 2024.
- 568 Ricky TQ Chen, Xuechen Li, Roger B Grosse, and David K Duvenaud. Isolating sources of dis-  
569 entanglement in variational autoencoders. *Advances in Neural Information Processing Systems*,  
570 2018.
- 571 Rewon Child. Very deep VAEs generalize autoregressive models and can outperform them on im-  
572 ages. In *International Conference on Learning Representations*, 2021.
- 573 Pierre Comon. Independent component analysis, a new concept? *Signal Processing*, 1994.
- 574 Bin Dai, Ziyu Wang, and David Wipf. The usual suspects? reassessing blame for VAE posterior  
575 collapse. In *International Conference on Machine Learning*, 2020.
- 576 Cian Eastwood and Christopher KI Williams. A framework for the quantitative evaluation of disen-  
577 tangled representations. In *International Conference on Learning Representations*, 2018.
- 578 Harrison Edwards and Amos Storkey. Censoring representations with an adversary. In *International  
579 Conference on Learning Representations*, 2016.
- 580 Babak Esmaeili, Hao Wu, Sarthak Jain, Alican Bozkurt, Narayanaswamy Siddharth, Brooks Paige,  
581 Dana H Brooks, Jennifer Dy, and Jan-Willem Meent. Structured disentangled representations. In  
582 *International Conference on Artificial Intelligence and Statistics*, 2019.
- 583 Chen Fang, Ye Xu, and Daniel N Rockmore. Unbiased metric learning: On the utilization of multiple  
584 datasets and web images for softening bias. In *International Conference on Computer Vision*,  
585 2013.
- 586 David Feldman, Avtar Singh, Jonathan L Schmid-Burgk, Rebecca J Carlson, Anja Mezger, An-  
587 thony J Garrity, Feng Zhang, and Paul C Blainey. Optical pooled screens in human cells. *Cell*,  
588 2019.
- 589  
590  
591  
592  
593

- 594 Hao Fu, Chunyuan Li, Xiaodong Liu, Jianfeng Gao, Asli Celikyilmaz, and Lawrence Carin. Cyclical  
595 annealing schedule: A simple approach to mitigating KL vanishing. *arXiv*, 2019.
- 596
- 597 Luke Funk, Kuan-Chung Su, Jimmy Ly, David Feldman, Avtar Singh, Britannia Moodie, Paul C  
598 Blainey, and Iain M Cheeseman. The phenotypic landscape of essential human genes. *Cell*, 2022.
- 599
- 600 Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François  
601 Laviolette, Mario March, and Victor Lempitsky. Domain-adversarial training of neural networks.  
602 *Journal of Machine Learning Research*, 2016.
- 603
- 604 Ishaan Gulrajani and David Lopez-Paz. In search of lost domain generalization. In *International  
605 Conference on Learning Representations*, 2021.
- 606
- 607 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-  
608 nition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*,  
609 2016.
- 610
- 611 Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv*, 2016.
- 612
- 613 Irina Higgins, Loic Matthey, Arka Pal, Christopher P Burgess, Xavier Glorot, Matthew M Botvinick,  
614 Shakir Mohamed, and Alexander Lerchner. beta-VAE: Learning basic visual concepts with a  
615 constrained variational framework. In *International Conference on Learning Representations*,  
616 2017.
- 617
- 618 Irina Higgins, David Amos, David Pfau, Sebastien Racaniere, Loic Matthey, Danilo Rezende, and  
619 Alexander Lerchner. Towards a definition of disentangled representations. *arXiv*, 2018.
- 620
- 621 Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected  
622 convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern  
623 Recognition*, 2017.
- 624
- 625 Aapo Hyvärinen and Erkki Oja. Independent component analysis: algorithms and applications.  
626 *Neural Networks*, 2000.
- 627
- 628 Michael I. Jordan, Zoubin Ghahramani, Tommi S. Jaakkola, and Lawrence K. Saul. An introduction  
629 to variational methods for graphical models. *Machine Learning*, 1999.
- 630
- 631 Christian Jutten and Jeanny Herault. Blind separation of sources, part i: An adaptive algorithm  
632 based on neuromimetic architecture. *Signal Processing*, 1991.
- 633
- 634 Ilyes Khemakhem, Diederik Kingma, Ricardo Monti, and Aapo Hyvarinen. Variational autoen-  
635 coders and nonlinear ICA: A unifying framework. In *International Conference on Artificial In-  
636 telligence and Statistics*, 2020.
- 637
- 638 Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron  
639 Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. In *Advances in Neural  
640 Information Processing Systems*, 2020.
- 641
- 642 Hyunjik Kim and Andriy Mnih. Disentangling by factorising. In *International Conference on  
643 Machine Learning*, 2018.
- 644
- 645 Diederik P. Kingma and Max Welling. Auto-encoding variational Bayes. In *International Confer-  
646 ence on Learning Representations*, 2014.
- 647
- 648 Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Bal-  
649 subramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, et al. Wilds: A  
650 benchmark of in-the-wild distribution shifts. In *International Conference on Machine Learning*,  
651 2021.
- 652
- 653 Lingjing Kong, Shaoan Xie, Weiran Yao, Yujia Zheng, Guangyi Chen, Petar Stojanov, Victor Akin-  
654 wande, and Kun Zhang. Partial disentanglement for domain adaptation. In *International Confer-  
655 ence on Machine Learning*, 2022.
- 656
- 657 Sébastien Lachapelle and Simon Lacoste-Julien. Partial disentanglement via mechanism sparsity.  
658 *arXiv*, 2022.

- 648 Sébastien Lachapelle, Divyat Mahajan, Ioannis Mitliagkas, and Simon Lacoste-Julien. Additive de-  
649 coders for latent variables identification and cartesian-product extrapolation. *Advances in Neural*  
650 *Information Processing Systems*, 2024.
- 651 Jeffrey T Leek, Robert B Scharpf, Héctor Corrada Bravo, David Simcha, Benjamin Langmead,  
652 W Evan Johnson, Donald Geman, Keith Baggerly, and Rafael A Irizarry. Tackling the widespread  
653 and critical impact of batch effects in high-throughput data. *Nature Reviews Genetics*, 2010.
- 654 Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M Hospedales. Deeper, broader and artier domain  
655 generalization. In *International Conference on Computer Vision*, 2017.
- 656 Romain Lopez, Jeffrey Regier, Michael I. Jordan, and Nir Yosef. Information constraints on auto-  
657 encoding variational bayes. In *Advances in Neural Information Processing Systems*, 2018.
- 658 Romain Lopez, Jan-Christian Huetter, Ehsan Hajiramezani, Jonathan K Pritchard, and Aviv Regev.  
659 Toward the identifiability of comparative deep generative models. In *Conference on Causal*  
660 *Learning and Reasoning*, 2024.
- 661 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Confer-*  
662 *ence on Learning Representations*, 2019.
- 663 Christos Louizos, Kevin Swersky, Yujia Li, Max Welling, and Richard S. Zemel. The variational  
664 fair autoencoder. In *International Conference on Learning Representations*, 2016.
- 665 Chaochao Lu, Yuhuai Wu, José Miguel Hernández-Lobato, and Bernhard Schölkopf. Invariant  
666 causal representation learning for out-of-distribution generalization. In *International Conference*  
667 *on Learning Representations*, 2021.
- 668 Taro Makino, Krzysztof Geras, and Kyunghyun Cho. Generative multitask learning mitigates target-  
669 causing confounding. *Advances in Neural Information Processing Systems*, 2022a.
- 670 Taro Makino, Stanislaw Jastrzebski, Witold Oleszkiewicz, Celin Chacko, Robin Ehrenpreis, Naziya  
671 Samreen, Chloe Chhor, Eric Kim, Jiyon Lee, Kristine Pysarenko, Beatriu Reig, Hildegard Toth,  
672 Divya Awal, Linda Du, Alice Kim, James Park, Daniel K. Sodickson, Laura Heacock, Linda Moy,  
673 Kyunghyun Cho, and Krzysztof J. Geras. Differences between human and machine perception in  
674 medical diagnosis. *Scientific Reports*, 2022b.
- 675 Sébastien Marcel and Yann Rodriguez. Torchvision the machine-vision package of torch. In *Inter-*  
676 *national Conference on Multimedia*, 2010.
- 677 Scott Mayer McKinney, Marcin Sieniek, Varun Godbole, Jonathan Godwin, Natasha Antropova,  
678 Hutan Ashrafian, Trevor Back, Mary Chesus, Greg S Corrado, Ara Darzi, et al. International  
679 evaluation of an AI system for breast cancer screening. *Nature*, 2020.
- 680 Krikamol Muandet, David Balduzzi, and Bernhard Schölkopf. Domain generalization via invariant  
681 feature representation. In *International Conference on Machine Learning*, 2013.
- 682 Tan H Nguyen, Dinkar Juyal, Jin Li, Aaditya Prakash, Shima Nofallah, Chintan Shah, Sai Chowdary  
683 Gullapally, Limin Yu, Michael Griffin, Anand Sampat, et al. Contrimix: Scalable stain color  
684 augmentation for domain generalization without domain labels in digital pathology. *arXiv*, 2023.
- 685 Jonas Peters, Peter Bühlmann, and Nicolai Meinshausen. Causal inference by using invariant pre-  
686 diction: identification and confidence intervals. *Journal of the Royal Statistical Society: Series B*  
687 *(Statistical Methodology)*, 2016.
- 688 Ali Razavi, Aäron van den Oord, Ben Poole, and Oriol Vinyals. Preventing posterior collapse with  
689 delta-vaes. *arXiv*, 2019.
- 690 Jennifer E Rood, Anna Hupalowska, and Aviv Regev. Toward a foundation model of causal cell and  
691 tissue biology with a perturbation cell and tissue atlas. *Cell*, 2024.
- 692 Andreas Ruepp, Brigitte Waegel, Martin Lechner, Barbara Brauner, Irmtraud Dunger-Kaltenbach,  
693 Gisela Fobo, Goar Frishman, Corinna Montrone, and H-Werner Mewes. CORUM: the compre-  
694 hensive resource of mammalian protein complexes—2009. *Nucleic Acids Research*, 2010.

- 702 David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-  
703 propagating errors. *Nature*, 1986.  
704
- 705 Shiori Sagawa, Pang Wei Koh, Tatsunori B. Hashimoto, and Percy Liang. Distributionally robust  
706 neural networks. In *International Conference on Learning Representations*, 2020.
- 707 Yuge Shi, Jeffrey Seely, Philip HS Torr, N Siddharth, Awni Hannun, Nicolas Usunier, and Gabriel  
708 Synnaeve. Gradient matching for domain generalization. In *International Conference on Learning  
709 Representations*, 2022.  
710
- 711 Baochen Sun and Kate Saenko. Deep coral: Correlation alignment for deep domain adaptation. In  
712 *ECCV 2016 Workshops*, 2016.
- 713 David Tellez, Geert Litjens, Péter Bándi, Wouter Bulten, John-Melle Bokhorst, Francesco Ciompi,  
714 and Jeroen Van Der Laak. Quantifying the effects of data augmentation and stain color normal-  
715 ization in convolutional neural networks for computational pathology. *Medical Image Analysis*,  
716 2019.
- 717 Damien Teney, Yong Lin, Seong Joon Oh, and Ehsan Abbasnejad. ID and OOD performance are  
718 sometimes inversely correlated on real-world datasets. *Advances in Neural Information Process-  
719 ing Systems*, 2024.  
720
- 721 Vladimir Vapnik. *The Nature of Statistical Learning Theory*. Springer, 1995.  
722
- 723 Julius Von Kügelgen, Yash Sharma, Luigi Gresele, Wieland Brendel, Bernhard Schölkopf, Michel  
724 Besserve, and Francesco Locatello. Self-supervised learning with data augmentations provably  
725 isolates content from style. *Advances in Neural Information Processing Systems*, 2021.
- 726 Yixin Wang, David Blei, and John P Cunningham. Posterior collapse and latent variable non-  
727 identifiability. *Advances in Neural Information Processing Systems*, 2021.  
728
- 729 Zitong Jerry Wang, Romain Lopez, Jan-Christian Hütter, Takamasa Kudo, Heming Yao, Philipp  
730 Hanslovsky, Burkhard Hoeckendorf, Rahul Moran, David Richmond, and Aviv Regev. Multi-  
731 contrastivevae disentangles perturbation effects in single cell images from optical pooled screens.  
732 *bioRxiv*, 2023.
- 733 Florian Wenzel, Andrea Dittadi, Peter Gehler, Carl-Johann Simon-Gabriel, Max Horn, Dominik  
734 Zietlow, David Kernert, Chris Russell, Thomas Brox, Bernt Schiele, et al. Assaying out-of-  
735 distribution generalization in transfer learning. *Advances in Neural Information Processing Sys-  
736 tems*, 2022.
- 737 Mitchell Wortsman, Gabriel Ilharco, Jong Wook Kim, Mike Li, Simon Kornblith, Rebecca Roelofs,  
738 Raphael Gontijo Lopes, Hannaneh Hajishirzi, Ali Farhadi, Hongseok Namkoong, et al. Robust  
739 fine-tuning of zero-shot models. In *Conference on Computer Vision and Pattern Recognition*,  
740 2022.  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

756 APPENDIX

757  
758 A EXPERIMENTS

759  
760 A.1 EXPERIMENTAL SETUP

761  
762 All of our experiments were done using a single NVIDIA A100 GPU on our institutions’ high-  
763 performance computing clusters.

764  
765 A.1.1 SUPERVISED CONTRASTIVE BLOCK DISENTANGLEMENT

766  
767 Our experimental setup for SCBD is remarkably similar across all of our experiments, which high-  
768 lights the generality of our approach. We resize the images to  $32 \times 32$  pixels and use a batch size of  
769 2,048. We set the temperature parameter in the supervised contrastive losses to  $\tau = 0.1$ . We adopt  
770 both of these practices from Khosla et al. (2020). For optimization, we use AdamW (Loshchilov &  
771 Hutter, 2019) with  $1 \times 10^{-4}$  learning rate and 0.01 weight decay. We chose these values because  
772 they resulted in stable training and validation curves across our experiments, and did not tune them  
773 extensively.

774 For domain generalization, we set  $D_{z_c} = D_{z_s} = 128$ , which we adopt from Khosla et al. (2020).  
775 We train for a maximum of 25,000 steps, and select the weights that minimize the validation loss.  
776 We obtain error bars by repeating each experiment with ten random seeds.

777 For batch correction, we set  $D_{z_c} = D_{z_s} = 64$  in order to ensure a fair comparison with the top-64  
778 PCA features of CellProfiler. To sample a minibatch, we first sample 256 distinct values of  $y$  from  
779 the class distribution of the training set, and then sample the same number of examples per value  
780 of  $y$ . This was necessary in order to ensure a large number of positive pairs with respect to  $y$  in  
781 our supervised contrastive losses, given that there are 5,050 classes. We trained for a maximum of  
782 150,000 steps, and evaluated on the test set using the weights that minimize the validation loss. We  
783 obtain error bars by repeating each experiment with three random seeds.

784 We use standard architectures such as ResNet-18 (He et al., 2016) and DenseNet-121 (Huang et al.,  
785 2017) for the encoders  $Enc_c(\mathbf{x})$  and  $Enc_s(\mathbf{x})$ . The projection networks  $Proj_c(\mathbf{r}_c)$  and  $Proj_s(\mathbf{r}_s)$   
786 are two-layer Multilayered Perceptrons (Rumelhart et al., 1986) with hidden sizes of  $D_{r_c}$  and  $D_{r_s}$ ,  
787 and GELU activations (Hendrycks & Gimpel, 2016). Our decoder  $Dec(\mathbf{z}_c, \mathbf{z}_s)$  architecture is shown  
788 in Appendix Table 2, with GELU activations (Hendrycks & Gimpel, 2016) between layers. We use  
789 an additive decoder (Lachapelle et al., 2024), and found this to be necessary to achieve sensible  
790 visualization results on CMNIST. That is, we define

$$791 \log p(\mathbf{x} \mid \mathbf{z}_c, \mathbf{z}_s) = Dec_c(\mathbf{z}_c) + Dec_s(\mathbf{z}_s),$$

792 where both  $Dec_c$  and  $Dec_s$  have the same architecture.

793  
794 Table 2: SCBD decoder architecture

795	Linear(64, 256 * (2 ** 2))
796	ConvTranspose2d(256, 256, 3, stride=2, padding=1, output_padding=1)
797	Conv2d(256, 256, 3, padding=1)
798	ConvTranspose2d(256, 256, 3, stride=2, padding=1, output_padding=1)
799	Conv2d(256, 256, 3, padding=1)
800	ConvTranspose2d(256, 256, 3, stride=2, padding=1, output_padding=1)
801	Conv2d(256, 256, 3, padding=1)
802	ConvTranspose2d(256, 256, 3, stride=2, padding=1, output_padding=1)
803	Conv2d(256, 128, 3, padding=1)
804	Conv2d(128, img_ch, 1)
805	

806  
807 A.1.2 VARIATIONAL AUTOENCODERS

808  
809 For our experiments with VAE-based approaches, we use the same experimental setup used in Wang  
et al. (2023), including the architecture and hyperparameters. We resize the images to  $64 \times 64$  pixels

and use a batch size of 1024. The encoder and decoder architectures are in Appendix Tables 3 and 4, with GELU activations (Hendrycks & Gimpel, 2016) between layers. Since the CMNIST images are  $32 \times 32$ , we modify the architectures to reduce the up- and down-sampling. For optimization, we use the AdamW (Loshchilov & Hutter, 2019) optimizer with  $1 \times 10^{-4}$  learning rate and 0.01 weight decay. We additionally skip gradients with a norm above  $1 \times 10^{12}$ , and clip gradients with a norm above  $1 \times 10^6$ , as done in Child (2021). We train for a maximum of 50,000 steps for domain generalization, and three epochs for batch correction, and select the weights with minimum validation loss. We report the validation and test performance across ten random seeds for domain generalization, and three random seeds for batch correction.

Table 3: VAE encoder architecture

---

```
Conv2d(img_c, 32, 3, stride=2, padding=1)
Conv2d(32, 32, 3, padding=1)
Conv2d(32, 64, 3, stride=2, padding=1)
Conv2d(64, 64, 3, padding=1)
Conv2d(64, 64, 3, stride=2, padding=1)
Linear(64 * (8 ** 2), 2 * 64)
```

---

Table 4: VAE decoder architecture

---

```
Linear(2 * 64, 64 * (8 ** 2))
ConvTranspose2d(64, 64, 3, stride=2, padding=1, out_padding=1)
Conv2d(64, 64, 3, padding=1)
ConvTranspose2d(64, 32, 3, stride=2, padding=1, out_padding=1)
Conv2d(32, 32, 3, padding=1)
ConvTranspose2d(32, img_c, 3, stride=2, padding=1, out_padding=1)
```

---

### A.1.3 OTHER BASELINES

Table 5: Hyperparameter search space

Condition	Hyperparameter	Search space
CMNIST	Learning rate	{0.0001, 0.001, 0.01}
	Weight decay	{0, 0.001, 0.01}
	Batch size	{32}
	Maximum epochs	{1, 20, 100}
Camelyon17-WILDS	Learning rate	{0.0001, 0.001, 0.01}
	Weight decay	{0, 0.001, 0.01}
	Batch size	{32}
	Maximum epochs	{5}
CORAL	Penalty weight	{0.1, 1, 10}
DANN	Penalty weight	{0.1, 1, 10}
IRM	Penalty weight	{1, 10, 100, 1000}
Fish	Pretrain steps	{1000, 10000}
	Meta learning rate	{0.001, 0.01, 0.1}
Group DRO	Step size	{0.01}

## A.2 DOMAIN GENERALIZATION

### A.2.1 COLORED MNIST

**Data generating process** The images are  $32 \times 32$  pixels, and are RGB. There are two training environments and a test environment. In the first training environment, which we label  $e = 0$ , we set the foreground pixels in the red channel to the value one, and those in the green and blue channels to the value  $y/|\mathcal{Y}|$ , where  $|\mathcal{Y}| = 10$  is the number of digits. For images with the digit zero,  $y = 0$ , so the digit is colored completely red. Then, as the digit increases from zero to nine, the digits are colored red, but with a decreasing intensity. In the second training environment, which we label  $e = 1$ , we set the foreground pixels in the green channel to the value one, and those in the red and blue channels to the value  $(|\mathcal{Y}| - 1 - y)/|\mathcal{Y}|$ . This has the effect of the digits being colored green, where the intensity increases with as the digit increases from zero to nine. In the test environment, the foreground pixels are set to one in all channels, which makes all of the digits white.

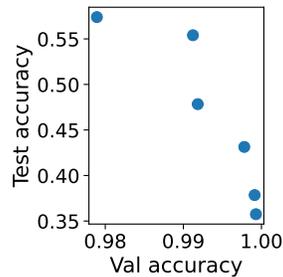


Figure 5: In- and out-of-distribution performance are negatively correlated on CMNIST, which satisfies the assumptions made by SCBD.

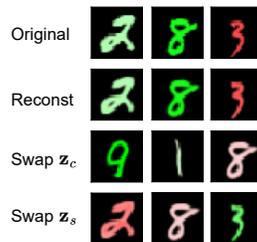


Figure 6: Counterfactual generation with iVAE. When we swap  $z_c$ , it changes the digit but not the color, and when we swap  $z_s$ , it changes the color but not the digit. iVAE generates better-looking images than SCBD, since the decoder is trained jointly with the encoder for iVAE.

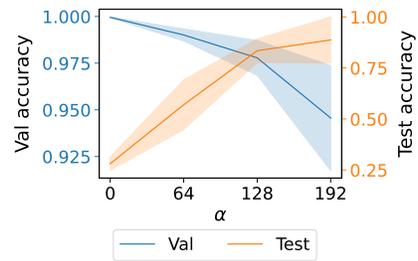


Figure 7: CMNIST results for an ablation in which we omit  $\mathbf{z}_s$ , and learn a single block of latent variables  $\mathbf{z}_c$  that are correlated with  $y$  and invariant to  $e$ . These results are similar to the model that learns  $\mathbf{z}_s$ . We use ResNet-18 encoders here, as we did in the main text.

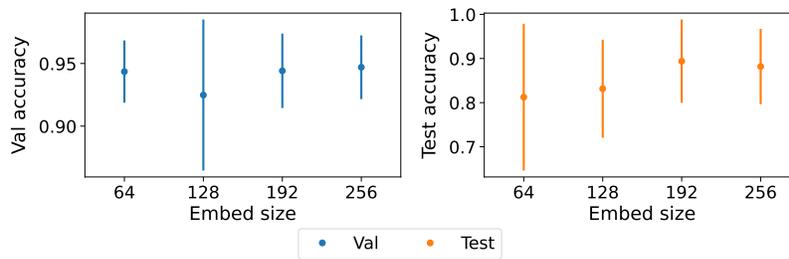


Figure 8: CMNIST embedding size ( $D_{\mathbf{z}_c}$  and  $D_{\mathbf{z}_s}$ ) ablation study for SCBD with ResNet-18 encoders and  $\alpha = 192$ . The results are relatively consistent across different embedding sizes.

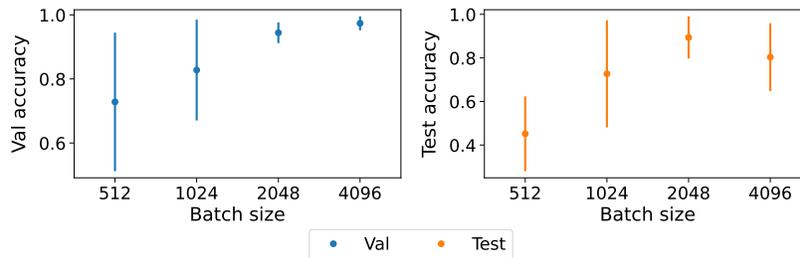


Figure 9: CMNIST batch size ablation study for SCBD with ResNet-18 encoders and  $\alpha = 192$ . The results are generally better for larger batch sizes, which was also observed by the authors of SCL (Khosla et al., 2020).

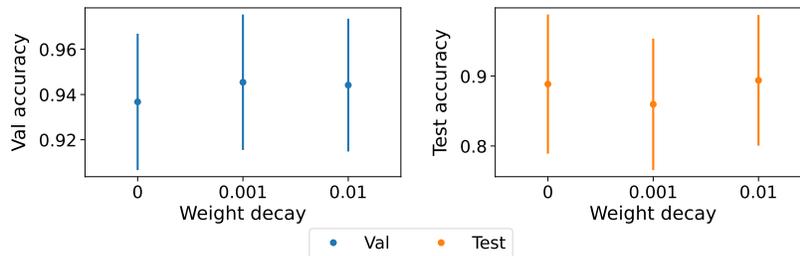


Figure 10: CMNIST weight decay ablation study for SCBD with ResNet-18 encoders and  $\alpha = 192$ . The results are relatively consistent across different degrees of weight decay.

## A.2.2 CAMELYON17-WILDS

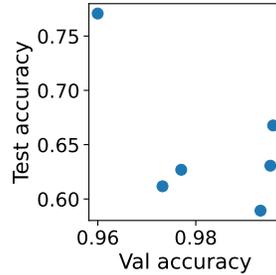


Figure 11: In- and out-of-distribution performance are negatively correlated on Camelyon17-WILDS. This is consistent with Teney et al. (2024), and therefore this dataset satisfies the assumptions made by SCBD.

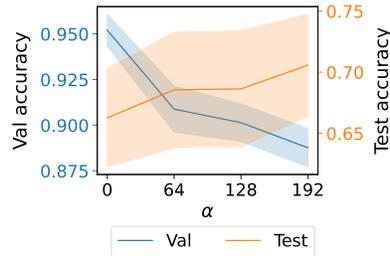


Figure 12: Camelyon17-WILDS results for SCBD using ResNet-18 encoders. The conclusions are the same as with the DenseNet-121 encoders.

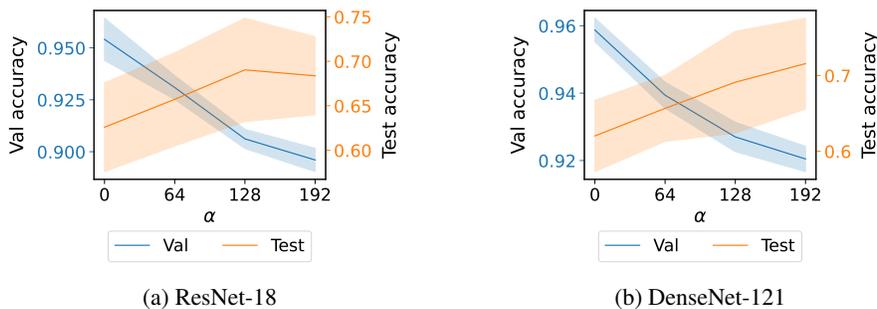


Figure 13: Camelyon17-WILDS results for an ablation in which we omit  $\mathbf{z}_s$ , and learn a single block of latent variables  $\mathbf{z}_c$  that are correlated with  $y$  and invariant to  $e$ . We observe a clean trade-off between validation and test accuracy with respect to  $\alpha$ , but the test accuracy error bars are larger than those of the model that includes  $\mathbf{z}_s$ .

1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038

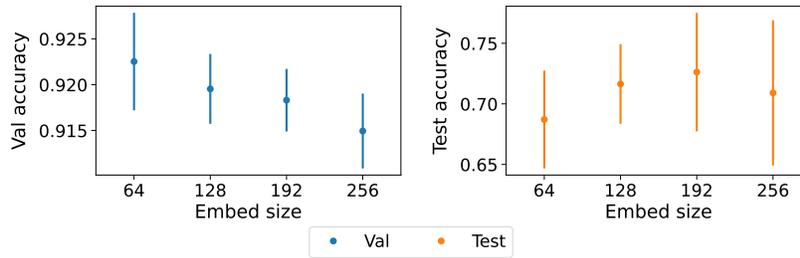


Figure 14: Camelyon17-WILDS embedding size ( $D_{z_c}$  and  $D_{z_s}$ ) ablation study for SCBD with DenseNet-121 encoders and  $\alpha = 192$ . The results are relatively consistent across different embedding sizes.

1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056

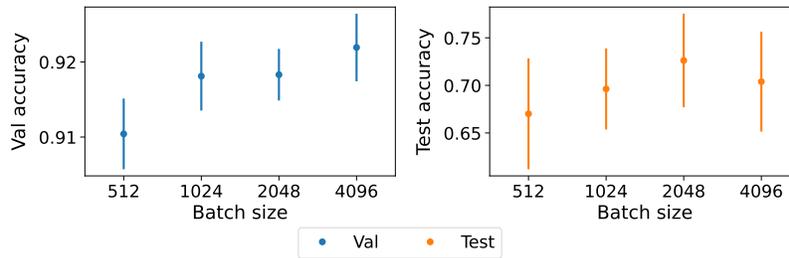


Figure 15: Camelyon17-WILDS batch size ablation study for SCBD with DenseNet-121 encoders and  $\alpha = 192$ . The results are generally better for larger batch sizes, which was also observed by the authors of SCL (Khosla et al., 2020).

1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079

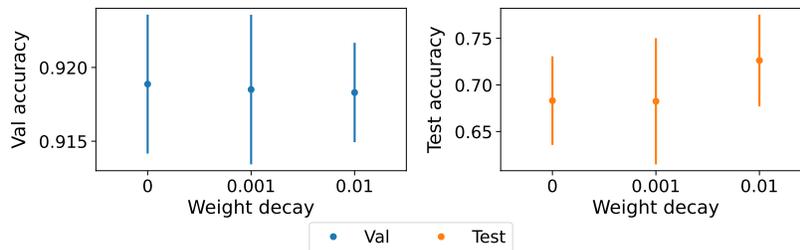
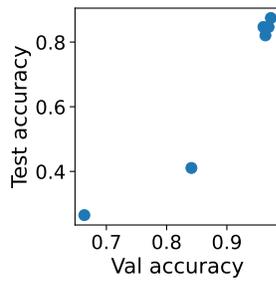
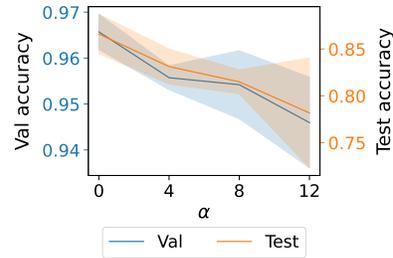


Figure 16: Camelyon17-WILDS weight decay ablation study for SCBD with DenseNet-121 encoders and  $\alpha = 192$ . The results are relatively consistent across different degrees of weight decay.

## A.2.3 PACS

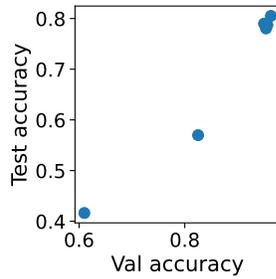


(a) In- and out-of-distribution performance are positively correlated, which violates the assumptions made by SCBD.

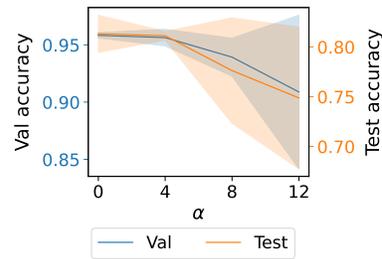


(b) Due to the violation of assumptions made by SCBD, increasing  $\alpha$  does not lead to a trade-off between in- and out-of-distribution performance.

Figure 17: PACS with art painting as the test domain.

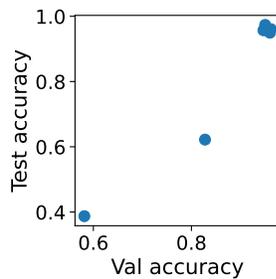


(a) In- and out-of-distribution performance are positively correlated, which violates the assumptions made by SCBD.

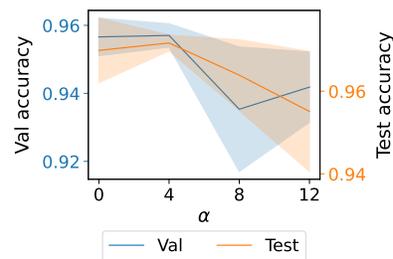


(b) Due to the violation of assumptions made by SCBD, increasing  $\alpha$  does not lead to a trade-off between in- and out-of-distribution performance.

Figure 18: PACS with cartoon as the test domain.



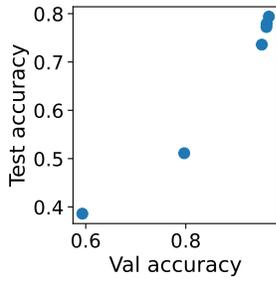
(a) In- and out-of-distribution performance are positively correlated, which violates the assumptions made by SCBD.



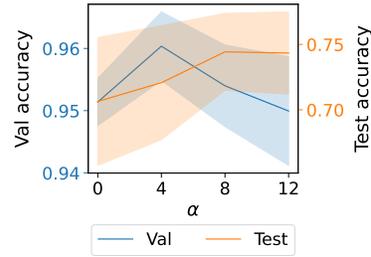
(b) Due to the violation of assumptions made by SCBD, increasing  $\alpha$  does not lead to a trade-off between in- and out-of-distribution performance.

Figure 19: PACS with photo as the test domain.

1134  
 1135  
 1136  
 1137  
 1138  
 1139  
 1140  
 1141  
 1142  
 1143  
 1144  
 1145  
 1146  
 1147  
 1148  
 1149  
 1150  
 1151  
 1152  
 1153  
 1154  
 1155  
 1156  
 1157  
 1158  
 1159  
 1160  
 1161  
 1162  
 1163  
 1164  
 1165  
 1166  
 1167  
 1168  
 1169  
 1170  
 1171  
 1172  
 1173  
 1174  
 1175  
 1176  
 1177  
 1178  
 1179  
 1180  
 1181  
 1182  
 1183  
 1184  
 1185  
 1186  
 1187



(a) In- and out-of-distribution performance are positively correlated, which violates the assumptions made by SCBD.



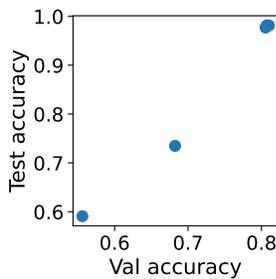
(b) Due to the violation of assumptions made by SCBD, increasing  $\alpha$  does not lead to a trade-off between in- and out-of-distribution performance.

Figure 20: PACS with sketch as the test domain.

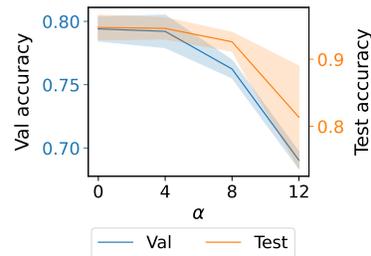
Table 6: Test accuracy (%) for PACS with three random seeds.

Algorithm	Art painting	Cartoon	Photo	Sketch	Average
SCBD ( $\alpha = 0$ )	$86.6 \pm 2.1$	$81.3 \pm 1.9$	$97.0 \pm 0.8$	$70.6 \pm 4.9$	83.9
ERM	$88.1 \pm 0.1$	$77.9 \pm 1.3$	$97.8 \pm 0.0$	$79.1 \pm 0.9$	85.7
CORAL	$87.7 \pm 0.6$	$79.2 \pm 1.1$	$97.6 \pm 0.0$	$79.4 \pm 0.7$	86.0
DANN	$85.9 \pm 0.5$	$79.9 \pm 1.4$	$97.6 \pm 0.2$	$75.2 \pm 2.8$	84.6
IRM	$85.0 \pm 1.6$	$77.6 \pm 0.9$	$96.7 \pm 0.3$	$78.5 \pm 2.6$	84.4
Group DRO	$86.4 \pm 0.3$	$79.9 \pm 0.8$	$98.0 \pm 0.3$	$72.1 \pm 0.7$	84.1

#### A.2.4 VLCS



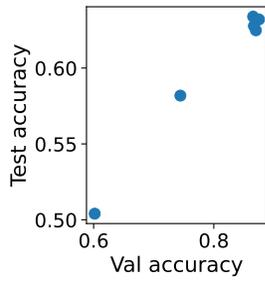
(a) In- and out-of-distribution performance are positively correlated, which violates the assumptions made by SCBD.



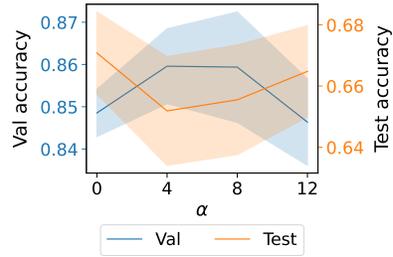
(b) Due to the violation of assumptions made by SCBD, increasing  $\alpha$  does not lead to a trade-off between in- and out-of-distribution performance.

Figure 21: VLCS with Caltech101 as the test domain.

1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197



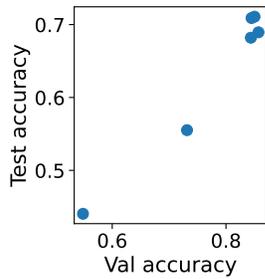
(a) In- and out-of-distribution performance are positively correlated, which violates the assumptions made by SCBD.



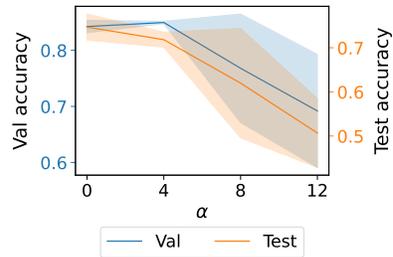
(b) Due to the violation of assumptions made by SCBD, increasing  $\alpha$  does not lead to a trade-off between in- and out-of-distribution performance.

Figure 22: VLCS with LabelMe as the test domain.

1202  
1203  
1204  
1205  
1206  
1207



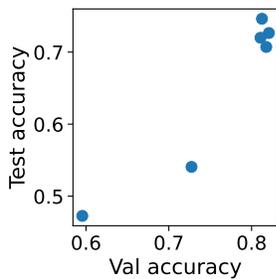
(a) In- and out-of-distribution performance are positively correlated, which violates the assumptions made by SCBD.



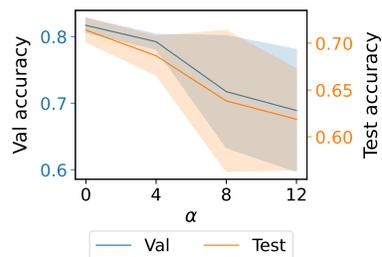
(b) Due to the violation of assumptions made by SCBD, increasing  $\alpha$  does not lead to a trade-off between in- and out-of-distribution performance.

Figure 23: VLCS with SUN09 as the test domain.

1222  
1223  
1224  
1225  
1226  
1227



(a) In- and out-of-distribution performance are positively correlated, which violates the assumptions made by SCBD.



(b) Due to the violation of assumptions made by SCBD, increasing  $\alpha$  does not lead to a trade-off between in- and out-of-distribution performance.

Figure 24: VLCS with VOC2007 as the test domain.

1241

Table 7: Test accuracy (%) for VLCS with three random seeds.

Algorithm	Caltech101	LabelMe	SUN09	VOC2007	Average
SCBD ( $\alpha = 0$ )	94.7 $\pm$ 1.8	67.1 $\pm$ 1.3	74.7 $\pm$ 3.0	71.4 $\pm$ 1.3	77.0
ERM	97.6 $\pm$ 1.0	63.3 $\pm$ 0.9	72.2 $\pm$ 0.5	76.4 $\pm$ 1.5	77.4
CORAL	98.8 $\pm$ 0.1	64.6 $\pm$ 0.8	71.7 $\pm$ 1.4	75.8 $\pm$ 0.4	77.7
DANN	98.5 $\pm$ 0.2	64.9 $\pm$ 1.1	73.1 $\pm$ 0.7	78.3 $\pm$ 0.3	78.7
IRM	97.6 $\pm$ 0.3	65.0 $\pm$ 0.9	72.2 $\pm$ 0.5	76.4 $\pm$ 1.5	78.1
Group DRO	97.7 $\pm$ 0.4	62.5 $\pm$ 1.1	70.1 $\pm$ 0.7	78.4 $\pm$ 0.9	77.2

### A.3 BATCH CORRECTION

**Batch prediction** We begin by computing the perturbation embedding for every example in the dataset, including the training, validation, and test sets. We then discard all examples for which we do not have a corresponding CellProfiler embedding. Using a randomly sampled 60% of the data as the training set, and the remaining data as the test set, we apply logistic regression to predict  $e$  given the embeddings, and report the F1 score on the test set.

**CORUM prediction** Our CORUM prediction task mirrors that of (Wang et al., 2023), with some modifications to ensure a fair comparison with CellProfiler. We begin by computing  $z_c$  for every single-cell image in the dataset, including the training, validation, and test sets. Then, we discard all embeddings for which we do not have a corresponding CellProfiler embedding. The median number of cells per gene is 6,000, and we want to average them to obtain a single embedding per gene. We have four sgRNA sequences for each perturbed gene, and 250 sgRNA sequences for the non-targeting control. We first average the  $z_c$ 's across cells for each sgRNA sequence, and then average the resulting sgRNA embeddings that correspond to the same gene. For each gene embedding, we subtract the non-targeting control embedding, then standardize such that each of the 64 components has mean zero and unit variance.

Then, we incorporate the CORUM database, which defines the pairs of genes that belong to the same protein complex. We discard all gene embeddings that are not in this database. We compute the cosine similarity between each pair of gene embeddings, and interpret it as the prediction that they belong to the same family. The prediction target is one if they belong to the same family according to the CORUM database, and zero otherwise. We turn the cosine similarities into binary predictions by across various prediction thresholds by using the  $i$ 'th percentile as the upper threshold and the  $100 - i$ 'th percentile as the lower threshold for each integer  $i \in \{80, \dots, 100\}$ . Finally, we use the binary predictions and prediction targets to obtain a precision and recall at each value of  $i$ , and plot the precision and recall curve.

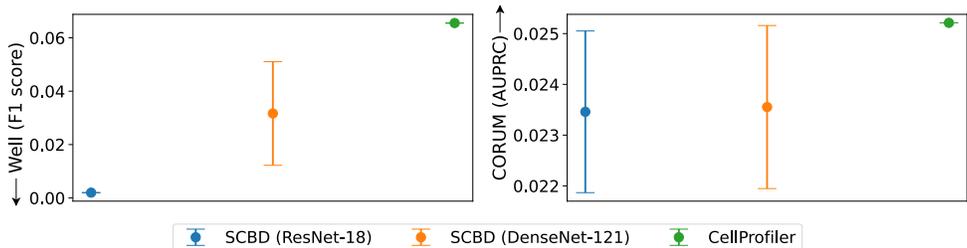


Figure 25: The left subfigure shows the performance of predicting the well label  $e$ , while the right subfigure represents the biological content. SCBD with DenseNet-121 encoders are less predictive of the well than CellProfiler, while retaining a similar level of biological information. However, the DenseNet-121 models are more sensitive to batch effects compared to the ResNet-18 models.

## 1296 B CODE

1297

1298

Here is our implementation of the supervised contrastive and invariance losses.

1299

1300

```
def supcon_loss(z, u, temperature):
    batch_size = len(z)
    u_col = u.unsqueeze(1)
    u_row = u.unsqueeze(0)
    mask_pos = (u_col == u_row).float()
    offdiag_mask = 1. - torch.eye(batch_size)
    mask_pos = mask_pos * offdiag_mask
    logits = torch.matmul(z, z.T) / temperature
    p = mask_pos / mask_pos.sum(dim=1, keepdim=True).clamp(min=1.)
    q = F.log_softmax(logits, dim=1)
    return F.cross_entropy(q, p)
```

1310

1311

```
def invariance_loss(zc, e, temperature):
    batch_size = len(zc)
    e_col = e.unsqueeze(1)
    e_row = e.unsqueeze(0)
    mask_pos = (e_col == e_row).float()
    mask_neg = 1. - mask_pos
    offdiag_mask = 1. - torch.eye(batch_size)
    mask_pos = mask_pos * offdiag_mask
    logits = torch.matmul(zc, zc.T) / temperature
    q = F.log_softmax(logits, dim=1)
    log_prob_pos = (q * mask_pos).mean(dim=1)
    log_prob_neg = (q * mask_neg).mean(dim=1)
    return (log_prob_pos - log_prob_neg).abs().mean()
```

1322

1323

1324

1325

1326

1327

1328

1329

1330

1331

1332

1333

1334

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349