# Domain Generalization - A Causal Perspective

#### A PREPRINT

Paras Sheth, Raha Moraffah, K. Selçuk Candan, Adrienne Raglin<sup>†</sup>, Huan Liu \*

November 8, 2022

### **ABSTRACT**

Machine learning models rely on various assumptions to attain high accuracy. One of the preliminary assumptions of these models is the independent and identical distribution, which suggests that the train and test data are sampled from the same distribution. However, this assumption seldom holds in the real world due to distribution shifts. As a result models that rely on this assumption exhibit poor generalization capabilities. Over the recent years, dedicated efforts have been made to improve the generalization capabilities of these models collectively known as – domain generalization methods. The primary idea behind these methods is to identify stable features or mechanisms that remain invariant across the different distributions. Many generalization approaches employ causal theories to describe invariance since causality and invariance are inextricably intertwined. However, current surveys deal with the causality-aware domain generalization methods on a very high-level. Furthermore, we argue that it is possible to categorize the methods based on how causality is leveraged in that method and in which part of the model pipeline is it used. To this end, we categorize the causal domain generalization methods into three categories, namely, (i) Invariance via Causal Data Augmentation methods which are applied during the data pre-processing stage, (ii) Invariance via Causal representation learning methods that are utilized during the representation learning stage, and (iii) Invariance via Transferring Causal mechanisms methods that are applied during the classification stage of the pipeline. Furthermore, this survey includes in-depth insights into benchmark datasets and code repositories for domain generalization methods. We conclude the survey with insights and discussions on future directions.

Keywords domain generalization · causality · vision · natural language processing · graphs

Machine learning (ML) models have achieved widespread success in variety of applications, including recommender systems [1], autonomous cars [2], and across various areas including, but not limited to, computer vision, graphs, and natural language processing. However, the success of these models is accompanied by various assumptions such as the independent and identical distributions or the i.i.d. assumption. According to this assumption the training and testing data are identically and independently distributed. In other words, the train and test data stem from the same distribution.

However, real-world data seldom abides by this assumption. Due to the dynamic nature of the systems that employ the machine learning models, the training data distribution might not be the same as the test data distribution and as a result the model's accuracy decreases. Under the machine learning paradigm, this phenomenon is known as *distribution shift*. For instance consider the task of digit classification as shown in Fig. 1. Given a model trained on a set of black and white handwritten digits [3] and evaluated on a set of colored handwritten digits [4] it is observed that the model performance drastically reduces, which raises the question *Why this happens?* It is well known that machine learning models leverage correlations among the different features to improve the prediction accuracy [5, 6]. However, in multiple real-world scenarios and the shown example, the model learns patterns (correlations) that are present in the training data but do not hold when evaluated against the test data. Thus, it results in the degrading performance.

One possible solution to address the distribution shift problem is to improve the *generalization* capabilities of the ML models which can be done in various ways including Domain Adaptation (DA) and Domain Generalization (DG).

<sup>\*</sup>Adrienne Raglin is with Army Research Laboratory (ARL). All other authors are with School of Computing and Augmented Intelligence, Arizona State University, Tempe, AZ, USA. (e-mail: {psheth5,rmoraffa,candan,huanliu}@asu.edu, adrienne.raglin2.civ@mail.mil)

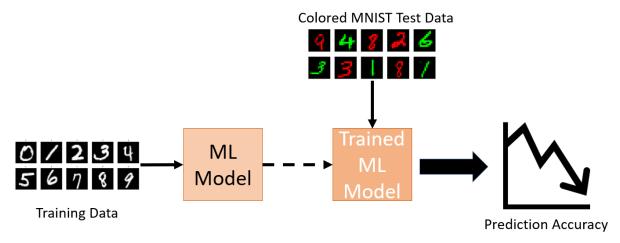


Figure 1: The task of digit classification from a domain generalization perspective. A model trained on MNIST dataset would fail to generalize when evaluated on ColoredMNIST, or RotatedMNIST.

Although domain adaptation methods lie outside the scope of this survey, we introduce them to distinguish between domain generalization and domain adaptation. Generally, both DA and DG focus on the same task, i.e., to improve the generalizability in ML models. However, the key difference between DA and DG methods is that the DA methods use a small number of unlabeled data points from the previously unseen test data (also known as target domain(s)) to fine-tune the model trained on the training data (also known as source domain(s)) [7, 8, 9]. However, in real-world settings the target data is usually not accessible beforehand to perform model adaptation. Thus, to deal with situations when the target domain data is unavailable, *domain generalization* methods were introduced [10].

With the distinction of DG and DA methods, we focus on understanding the working mechanism of DG methods, and how causality aids in the process. To understand this, let's consider the aforementioned digit classification scenario. Although the train domain data consisted of black and white images of digits and the target domain consisted of colored digits, as humans we are still able to classify the digits even after the distribution shift. Even though the digits are colored in the target domain, the shape of the digit, which is crucial in determining the digit, remains invariant. This implies that even in the presence of distribution shifts there exists a set of features that are invariant and are crucial for the prediction performance. Furthermore, it is well established that causality and invariance are tightly linked to each other, i.e., one of the dimensions of causality is invariance [11, 12]. Thus, causality can be a useful tool in capturing the invariance present in the data, justifying the range of methods that have leveraged different causal theories for improving the generalization capabilities of models [4, 13].

Since there exist multiple causality aware domain generalization methods, the next step should be to understand how these methods leverage causality and how they differ among themselves. Although majority of the existing surveys discuss causal domain generalization methods [14, 15, 16] they mostly club all the methods under the same umbrella term, i.e., causality aware domain generalization methods or causal representation methods. We argue that these methods can be better classified based on where the causal theories are leveraged during the entire model pipeline. To this end, we categorize these papers broadly into three different categories: (i) Invariance via Causal Data Augmentation methods which are applied during the data pre-processing stage, (ii) Invariance via Causal representation learning methods that are utilized during the representation learning stage, and (iii) Invariance via Transferring Causal mechanisms methods that are applied during the classification stage of the pipeline. Furthermore, Fig. 2 shows the detailed breakdown of these categories into their corresponding subcategories.

The remaining part of the survey is categorized as follows. Section 1 covers the problem definition and causal preliminaries, Section 2 covers the different categories and the sub-categories helping the reader understand how causality can be leveraged in different parts of the model pipeline and how each category tackles its corresponding challenges. Section 3 covers the Causal DG methods developed for the graphs and Natural Language Processing (NLP). Section ?? covers the benchmark datasets and evaluation schemes employed by various methods for evaluating the performances. This section also covers the publicly available code repositories which could aid researchers in developing and testing their own models. We finally conclude the survey with Section ??.

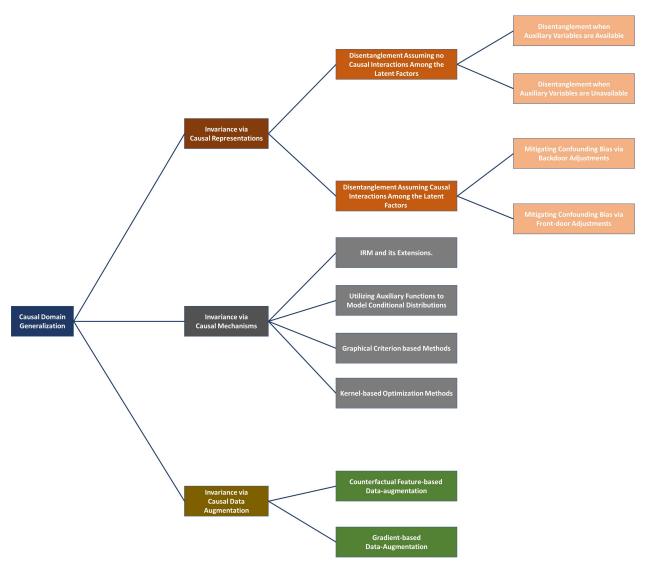


Figure 2: The categorization of Causal Domain Generalization techniques.

### 1 Problem Definition and Preliminaries

### 1.1 Problem Definition

Consider X as the set of features, Y as the set of labels, and D as the set of domain(s) with sample spaces  $\mathcal{X}$ ,  $\mathcal{Y}$ , and  $\mathcal{D}$ , respectively. A domain is defined as a joint distribution  $P_{X,Y}$  on  $\mathcal{X} \times \mathcal{Y}$ . Let  $P_X$  represent X's marginal distribution,  $P_{X|Y}$  represent the class-conditional distribution of X given Y, and  $P_{Y|X}$  represent the posterior distribution of Y given Y.

A domain generalization model's purpose is to learn a predictive model  $f: \mathcal{X} \to \mathcal{Y}$ . However, while dealing with domain generalization, the common assumption implies that training data was obtained from a finite subset of the possible domains  $D_{train} \subset \mathcal{D}$ . Furthermore, the number of training domains is given by K, and  $D_{train} = \{d_i\}_{i=1}^K \subset \mathcal{D}$ . As a result, the training data is sampled from a distribution  $P[X,Y\mid D=d_i] \quad \forall i\in [k]$ . The domain generalization model, then aims at utilizing only source (train) domain(s) data with the goal of minimizing the prediction error on a previously unseen target (test) domain. The corresponding joint distribution of the target domain  $D_{test}$  is given by  $P_{X,Y}^{D_{test}}$  and  $P_{X,Y}^{D_{test}} \neq P_{XY}^{(k)}, \forall k \in \{1, \dots, K\}$ . Ideally, the goal is to learn a classifier that is optimal for all domains  $\mathcal{D}$ .

#### 1.2 Preliminaries

Most causality-aware domain generalization models aim to mitigate the reliance on spurious correlations by accounting for the confounding effects. For instance, some works mitigate the confounding bias induced by the non-causal features on the causal features and the labels. This leads to better generalizability of the machine learning model. Confounding variables refer to those sets of variables that may be observed or unobserved, directly influencing the supposed cause and effect variables.

The confounding bias can be accounted for with either backdoor adjustment or front-door adjustment based on the nature of the problem. Therefore, we begin by defining the backdoor criterion.

**Definition 1.** Given a Directed Acyclic Graph (DAG) G, a set of nodes Z, and a pair of nodes, namely, X and Y, we say that Z satisfies the backdoor criterion relative to X and Y if:

- no node in Z is a descendant of Y
- Every path between X and Y that contains an arrow in Y is blocked by Z.

If Z satisfies the backdoor criterion for X and Y, the causal effect between X and Y is identifiable. The causal effect from X to Y can be formulated as,

$$P(Y \mid X) = \sum_{z} P(Y \mid X, Z)P(Z). \tag{1}$$

Similarly, another way to account for the confounding variables is the front-door criterion.

**Definition 2.** Given a Directed Acyclic Graph (DAG) G, a set of nodes Z, and a pair of nodes, namely, X and Y, we say that Z satisfies the front-door criterion relative to X and Y if:

- All directed paths from X to Y are intercepted by Z;
- There exists no unblocked backdoor path from X to Z; and
- X blocks all possible backdoor paths from Z to Y.

As per the front-door adjustment, if Z satisfies the front-door criterion for X and Y, and if P(X, Z) > 0, then the causal effect between X and Y is identifiable. The causal effect from X to Y can be formulated as,

$$P(Y \mid do(X)) = \sum_{z} \sum_{x'} P(Y \mid Z, X') P(X') P(Z \mid X).$$
 (2)

### 2 Domain Generalization via Invariance Learning

Causal Domain Generalization methods aim to leverage causal theories to improve the generalization capabilities of models. The causal theories are utilized in various stages of the standard machine learning pipeline. To this end, we categorize the Causal DG methods based on how and where the causality aspects are utilized. To this end, we propose three categories, namely, (1) Invariance via Causal Data Augmentation methods which are applied during the data pre-processing stage, (2) Invariance via Causal representation learning methods that are utilized during the representation learning stage, and (3) Invariance via Transferring Causal mechanisms methods that are applied during the classification stage of the pipeline.

Furthermore, we present the different subcategories associated with each category (when applicable) and discuss the methods for each corresponding sub-category. A summary of the categorization can be seen in Fig. 2.In the following sub-sections, we provide a comprehensive and detailed review of these methods corresponding to the above order and discuss their differences and theoretical connections.

These frameworks learn generalizable feature representations across different domains. Often, these representations are interpreted as causal feature representations. One possible way to learn such causal representations is to disentangle the input representations into causal and non-causal feature representations. By doing so, the domain shift can be justified as interventions on the non-causal features. In such a condition, the primary goal is to minimize a loss that is robust under changes in the distribution of these style (non-causal) features. For example, in Fig.3, the learned input representations can be divided into causal and non-causal features. Utilizing only the causal features for the task can improve generalizability.

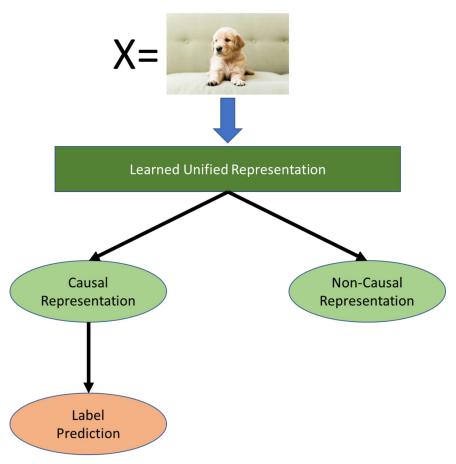


Figure 3: A simple causal graph that disentangles the input representations into causal and non-causal representations. Only the causal representations are used for prediction.

#### 2.1 Invariance via Causal Data Augmentation

In this section, we present the frameworks that achieve invariance via causal data augmentation. Frameworks in this category utilize causal features and augment the data by considering all possible confounders/ spurious variables. Although these methods eventually aim to learn causal representations, the mechanisms they employ (such as identifying the features to augment on) follow a causal procedure thus promoting such methods to have their own category.

### 2.1.1 Counterfactual Feature based Data Augmentation

Zhang et al. [17] propose a generalizable framework for the task of human pose estimation. The proposed framework aims to leverage generation of counterfactuals by intervening on the domain variable. In particular, these counterfactuals are generated by changing the domain variable. This framework enforces three main criteria to learn the causal representations: (1) Producing low-error prediction over the observed samples; (2) Producing low-error prediction over counterfactual samples; and (3) The counterfactual and observed feature representations are similar. The authors propose to implement the framework in two branches, namely observed branch and counterfactual branch. In the observed branch, images from source domain are fed to a feature extractor (encoder) to get the observed feature representation distribution. In the counterfactual branch, a GAN-based architecture is utilized to learn the distribution of counterfactuals from a ground-truth pose and random noise. Both observed and counterfactual representations are then fed to predictors to ensure they have high predictive power. To minimize the distance between observed and counterfactual representations a smooth- $l_1$  distance is utilized. The objective of this framework is given below:

$$\min_{\theta_f, \theta_h} \mathbb{E}_{(x,y,u) \sim (p(x), p(y), p(u)} \mathcal{L}_F(h(f(x)), y) 
\lambda_1 \mathcal{L}_{CF}(h(g(u, y)), y) + \lambda_2 \mathcal{L}_{dist}(f(x), g(u, y)),$$
(3)

where  $\mathcal{L}_F$  and  $\mathcal{L}_{CF}$  denote the prediction loss over observed and counterfactual representations,  $\lambda_1$  and  $\lambda_2$  are hyperparameters.

Chen et al. [18] propose the Interventional Emotion Recognition Network (IERN) to improve the visual emotion recognition framework's generalization via causal data augmentation. IERN first disentangles the input image into the context features (confounding features) and emotion features (causal features). Finally, the framework proposes to debias the classifier via backdoor criterion using the following objective:

$$\mathcal{L}_{cl} = \min_{f_c^{\theta}, g_e^{\theta}, f_b^{\theta}} (l_{CE}(\frac{1}{N_c}(\sum_{i=1}^{N_C} f_c(g_r(g_e(f_b(x)), C_i)), y_e))), \tag{4}$$

where  $f_c$  denotes the classifier,  $g_e(f_b(x))$  represent the emotion features,  $C_i$  are the confounders, ad  $N_C$  is the number of confounders. Ouyang et al. [19], propose a causal data augmentation approach to single-source domain generalization problem, where training data is only available from one source domain. The authors propose that the performance deterioration under domain shift may arise due to shifted domain-dependent features or shifted-correlation effect which is induced due to the presence of confounders. To deal with this problem, the authors propose to utilize causal intervention to augment the data and improve the robustness and generalization of the model. The proposed model consist of two parts: (1) global intensity non-linear augmentation (GIN) technique that utilizes randomly-weighted shallow convolutional networks to transforms images while keeping the shapes invariant; (2) interventional pseudocorrelation augmentation (IPA) technique that removes the confounder via re-sampling appearances of confounded objects independently. To augment the data, the framework first applies GIN on the input images to get new appearances. Then the two GIN-augmented of the same image are blended via IPA in a spatially-variable manner. Mitrovic et al. [20] propose a novel self-supervised objective, Representation Learning via Invariant Causal Mechanisms (RELIC) based on a causal analysis of contranstive learning frameworks. To do so, the authors propose to model the data generation process similar to the causal graph presented in Fig.3. The graph indicates that the data is generated from content (causal) and style (non-causal) variables and that only the only content (causal) is informative of the downstream tasks (label prediction). Given the causal graph, the paper indicates that the optimal representation should be invariant predictor of proxy targets on correlated, not causally related features. Since none of the causal or correlational variables are known, the authors utilize data augmentations to simulate interventions on the styles (correlational) variables. Finally, the paper proposes a regularizer to enforce the invariance under data augmentations as follows:

$$\min \mathbb{E}_{X \sim p(X)} \mathbb{E}_{a_{lk}, a_{qt} \sim \mathcal{A} \times \mathcal{A}} \sum_{b \in \{a_{lk}, a_{qt}\}} \mathcal{L}_b(Y^R, f(X))$$

$$\text{s.t.} KL(p^{do(a_{lk})}(Y^R \mid f(X)), p^{do(a_{qt})}(Y^R \mid f(X))) \leq \rho,$$

$$(5)$$

where  $A = \{a_1, ..., a_m\}$  is a set of data augmentations generated by intervening on the style variables,  $\mathcal{L}$  is the proxy task loss and KL is the Kullback-Leibler (KL) divergence. This KL-divergence based regularizer enforces that the prediction of the proxy targets is invariant across data augmentations.

### 2.1.2 Gradient-based Data Augmentation

Bai et al. [21] argue that many domain generalization methods work well for one dataset but perform poorly for others. They claim this happens because the domain generalization problems have multiple dimensions - correlation shift and diversity shift. Correlation shift is when the labels and the environments are correlated, and the relations change across different environments. Diversity shift means the data comes from different domains, thus having significantly different styles. For instance, the sketch of a horse, a cartoon horse, an image of a horse, and an art of a horse all represent horses in different styles. Furthermore, real-world data exists, which is a mixture of these shifts. To handle these dimensions simultaneously, they propose a novel decomposed feature representation and semantic augmentation approach. First, the proposed method decomposes the representations of the input image into context and category features. Then, they perform gradient-based semantic augmentation on context features, representing attributes, styles, and more; to disentangle the spurious correlation between features. The semantic data augmentation is performed by adversarially perturbing the feature space of the context related features of the original sample as follows:

$$z_{i}^{c} = z_{i}^{c} + \alpha_{i} \cdot \epsilon \cdot \frac{\nabla_{z_{i}^{c}}(l(h_{\theta_{c}})(z_{i}^{c}, c_{i})))}{\|\nabla_{z_{i}^{c}}(l(h_{\theta_{c}})(z_{i}^{c}, c_{i})))\|},$$
(6)

where  $z_i^c$  is the context feature representation,  $h_{\theta_c}$  is the context feature discriminator,  $\epsilon$  is a hyperparameter which controls the maximum length of the augmentation vectors, and  $\alpha_i$  is randomly sampled from [0,1]. Unlike methods in the earlier category, this work does not aim to generate counterfactuals to improve generalization, rather they perform gradient based augmentation on disentangled context features to eliminate distribution shifts for various generalization tasks.

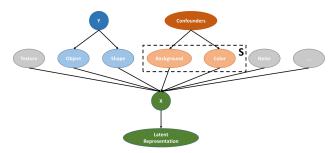


Figure 4: The data generating process as shown in [22]. As per the authors, the input X is generated by a collection of unobserved factors (such as texture, object, shape, and so on). Only Y, X, and possibly the confounders are observed during training time. Y represents the variable that needs to be predicted. The generative factors are assumed to not have any causal relation among themselves. This work also assumes that the label and the confounders have an effect on the generative factors. S represents the style variables.

### 2.2 Invariance via Learning Causal Representations

Causal DG methods aim to capture invariance with the aid of causal theories. Majority of these methods aim at learning a causal representation of the data that is highly predictive of the downstream task. For instance, consider the task of digit classification as discussed in the Introduction. Despite applying various transformations to the digits, (such as rotating the digits, coloring the digits, and so on) humans are still able to correctly associate digits with the labels. This phenomenon could be attributed to the fact that humans utilize the features that do not change during these transformations such as the shape of the digit. Thus, these features can be considered causal. However, learning these representations is not straightforward as these features are not directly observable making it challenging to guide the machine. One possible way to guide the model to rely on causal representations is with the aid of disentanglement. The range of works that consider disentanglement can be further divided into two parts namely, works that consider disentanglement with no causal interactions among the latent factors, and works that consider disentanglement with causal interactions among the latent factors. In this section we discuss these different categories.

### 2.2.1 Disentanglement Assuming No Causal Interactions Among the Latent Factors

This section discusses the works that assume the latent factors (i.e. the causal and non-causal factors) bare no causal interactions amongst each other. Majority of the papers utilize a simple causal graph as seen in Fig. 3. The input features are divided into causal and non-causal features.

### Disentanglement when Auxiliary Variables are Available

A range of variables utilize auxiliary variables to aid the disentanglement process. For instance, in the task of image classification, auxiliary variables could be additional cues about the object, or they could be background variables indicating the image background. Thus, auxiliary variables can aid in distinguishing causal and non-causal features. The authors in [13, 23, 24] aim to utilize auxiliary variables to separate causal from non-causal features, and learn the representations accordingly. For example, in grouped observations for certain datasets the same object (ID) is seen under multiple situations [23], which can guide the prediction to be based more on the latent core (causal) characteristics and less on the latent style (non-causal) features by penalizing between-object variance of the prediction less than variation for the same object. The authors of [23] contend that direct interventions on style elements frequently result in the creation of a new domain. So, if  $F_0$  represents the joint distribution of the  $(ID, Y, X^{style})$  in the training distribution, then intervening on  $X^{style}$  yields a new joint distribution of the  $(ID, Y, \tilde{X}^{style})$  indicated by F. As a result, we obtain the following class of distributions:

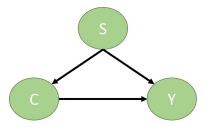
$$\mathcal{F}_{\mathcal{E}} = \{ F : D_{\text{style}} (F_0, F) \le \xi \} \tag{7}$$

where  $D_{\text{style}}(F_0, F)$  is the distance between the two distributions. The primary goal is to optimize a worst-case loss over this distribution class. This loss can be formulated as,

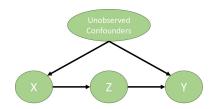
$$L_{\xi}(\theta) = \sup_{F \in \mathcal{F}_{\xi}} E_F \left[ \ell\left(Y, f_{\theta}(X)\right) \right] \tag{8}$$

For arbitrary strong interventions on the style features, the loss would be given by,

$$L_{\infty}(\theta) = \lim_{\xi \to \infty} \sup_{F \in \mathcal{F}_{\xi}} E_F \left[ \ell \left( Y, f_{\theta}(X) \right) \right] \tag{9}$$



(a) Causal graph for Backdoor adjustment.



(b) Causal graph for Front-door adjustment.

Minimizing this loss guarantees an accurate prediction that performs well even for significant shifts in the conditional distribution of style features. Rather than pooling over all examples, CoRe [23] exploits the ID variable to penalize the loss function. The overall objective function is given by,

$$\hat{\theta}^{\text{core}}(\lambda) = \operatorname{argmin}_{\theta} \hat{E}\left[\ell\left(Y, f_{\theta}(X)\right)\right] + \lambda \cdot \hat{C}_{\theta} \tag{10}$$

where  $\hat{C}_{\theta}$  is a conditional variance penalty of the form

$$\hat{C}_{f,\nu,\theta} := \hat{E}[\widehat{\operatorname{Var}}(f_{\theta}(X) \mid Y, \operatorname{ID})^{\nu}]$$
(11)

for the conditional-variance-of-prediction, and

$$\hat{C}_{\ell,\nu,\theta} := \hat{E} \left[ \widehat{\operatorname{Var}} \left( \ell \left( Y, f_{\theta}(X) \right) \mid Y, \operatorname{ID} \right)^{\nu} \right]$$
(12)

for the conditional-variance-of-loss.  $f_{\theta}(X)$  is the representation of the input X, Y is the image label, ID is the identifier label or the object label, and  $\nu \in \{1/2, 1\}$ . Mahajan et al. [13] argue that for representations, the class-conditional domain invariant objective is inadequate. They add to the CoRe framework by including an approach for when objects (ID variables) are not detected. When the stable feature distribution differs across domains, the class-conditional aim is insufficient to learn the stable features. To solve this issue, they use a causal graph to express within-class variance in stable features. The authors propose using the causal graph to learn the representation regardless of the classification loss. Liu et al. [25] claim that CoRe adds no new generative modeling efforts at the expense of limited capability for invariant causal processes They further contend that deep learning methods fail to generalize to unexplored domains because the representations they learnt blend semantic and stylistic information owing to false correlations. They propose using a causal generative model that follows a causal approach to describe the semantic and stylistic aspects independently to address these issues. Similar to CoRe, Makar et al. [26] suggest that identifying invariant characteristics is a tough undertaking since it is impossible to distinguish the effect of non-causal factors without extra supervision. Instead of ID labels, the authors recommend employing auxiliary labels (such as picture background labels) to provide information about the irrelevant component that is available during training but not during testing. The authors offer a method for using these auxiliary labels to build a predictor whose risk is roughly invariant over a well-defined range of test distributions.

Neet et al. [27] argue that causality-aware domain generalization frameworks impose regularization constraints to learn the invariance. However, when datasets with multi-attribute shifts are considered, deciding which regularization could lead to learning the invariance is difficult. To overcome this problem, the authors propose *Causally Adaptive Constraint Minimization* (CACM). CACM, leverage the information provided by multiple independent shifts across attributes, assuming structural knowledge of the shifts. By identifying the correct constraints, CACM applies them as regularizers in the overall objective function. Moreover, the attributes correlate with the label that can change between train and test data. The model aims to learn a risk-invariant predictor that obtains minimum risk on all distributions. Mathematically,

$$g_{\text{rinv}} \in \arg\min_{g \in G_{rinv}} R_P(g) \forall P \in \mathcal{P}$$
where  $G_{rinv} = \{g : R_P(g) = R_{P'}(g) \forall P, P' \in \mathcal{P}\},$ 

$$(13)$$

where R is the risk of predictor g, and  $\mathcal{P}$  represents all the distributions. The authors divide the auxiliary attributes A into three types of attributes  $A_{ind}$ , which represents the attributes that are correlated with the label,  $A_{ind}$ , the attributes that are independent of the label and E denotes the domain. The shifts are based on the relationship between the auxiliary attributes A and the label Y. Based on these distinctions of attributes, the authors define four kinds of shifts that are possible based on the causal graph. They are Independent, Causal, Confounded, and Selected. First, the authors identify the conditional independence constraints satisfied by the causal features and enforce that learned representation  $\phi$  should follow the same constraint. The authors leverage the d-seperation [28] strategy and first, for every observed variable  $V \in \mathcal{V}$  in the graph, check whether  $(X_c, V)$  are d-separated. If not, check whether  $(X_c, V)$  are d-separated

conditioned on any subset of the remaining observed variables in  $V \setminus \{V\}$ . Finally, the model applies those constraints as a regularizer to the standard ERM loss. Thus, the final objective function is formulated as,

$$g_1, \phi = \underset{g_1, \phi}{\operatorname{arg\,min}}; \quad \ell\left(g_1(\phi(\boldsymbol{x})), y\right) + \lambda^*(\text{ RegPenalty }),$$
 (14)

where  $\ell$  is the classification loss,  $\lambda$  is a hyperparameter, and RegPenalty is the regularization penalty. The penalty depends on the type of distribution shift for each attribute.

#### Disentanglement when Auxiliary Variables are Unavailable

Although auxiliary variables can aid with causal disentanglement, these variables are not always easily available. A series of work aim to tackle the causal disentanglement problem in the absence of auxiliary variables [22, 27, 29, 30]. Chevalley et al. [22] argue that the invariant representations are reformulated as a feature of a causal process, and propose a regularizer that guarantees invariance via distribution matching. The proposed framework, in particular, describes the underlying generation process using the DAG shown in Fig. 4. Different domains are then represented in the provided DAG via soft-intervention on the domain variable. The invariant representations are therefore characterized as those that have no complete causal influence on the domain variable. The classification is formulated as:

$$\min_{Z=f(X)} \mathcal{L}(Y, c(Z))$$
s.t.  $\mathbb{E}_{N_d, N_d'}[\operatorname{dist}(p^{do(d=N_d)(Z), p^{do(d=N_d')(Z)}})],$ 

$$(15)$$

where Z is the learned representation and dist denotes the distance between the distributions and  $N_d$ ,  $N_d'$  are the interventions on the domain variable. One way to disentangle causal and non-causal features in the presence of multiple domains, but the absence of auxiliary variables would be to utilize contrastive learning. The primary assumption under this setting is that the non-causal feature representations are similar for instances from the same domain. Thus, by guiding the machine learning model to learn non-causal representations, we can learn causal representations by learning orthogonal representations to the non-causal representations. In this setting, the objective function is usually represented as.

$$\mathcal{L} = \mathcal{L}^{cls} + \mathcal{L}^{con} \tag{16}$$

where  $\mathcal{L}^{cls}$  represent the classification loss, and  $\mathcal{L}^{con}$  represents the contrastive loss.  $\mathcal{L}^{cls}$  is formulated similar to Eq. 8 as it aims to predict the image label using the representation of the causal factors.  $\mathcal{L}^{con}$  is formulated as,

$$\mathcal{L}_{i,j}^{con} = -\log \frac{\exp(\sin(z_i, z_j)/\tau)}{\sum_{k=1}^{2N} \exp(\sin(z_i, z_k)/\tau)}$$
(17)

where  $\tau$  is the temperature normalization factor, and sim is the similarity function. The intuition here is that we want the similar representations  $z_i$  and  $z_j$  to be close to each other, and the dissimilar representations  $z_i$  and  $z_k$  to be more distant from each other. Recent works, including [24, 29] leverage this assumption to identify the causal features. Chen et al. [29] Concentrate on utilizing both semantic (causal) and stylistic (non-causal) characteristics. This paper's primary hypothesis is that instances from the same domain should exchange style information. This aids in the efficient disentanglement of style features, easing the search for actual semantic features with a low degree of freedom. Following the acquisition of style data, the network learns semantic information in orthogonal directions to the domain style. As a result, the network avoids overfitting style cues and instead concentrates on learning semantic elements. On the other hand, Trivedi et al. [24] claim that standard domain generalization methods are inapplicable to computer games because games are more than just graphics, gaming images include functional qualities related with the game genre as well as aesthetic features unique to each game. As a result, representations derived from a pre-trained model perform well on the game on which it was trained, resulting in poor generalization. To overcome this issue, the authors recommend using contrastive learning. The authors use the causal graph to distinguish between content and stylistic aspects, stating that the game genre is defined by the content elements.

The ways mentioned earlier are helpful for the traditional machine learning setting. However, in an online learning scenario, where an agent continually learns as it interacts with the world, it is more challenging to identify the causal features. To tackle this problem, Javed et al. [30] propose a framework for detecting and removing spurious features on a continuous basis. The framework's basic premise is that the association between a false characteristic and the goal is not continuous over time. The authors also contend that in order to infer a causal model from observable data, an agent must go through three phases. First, the agent must understand the basic data representations. Second, the agent must deduce the causal relationships between the variables. Finally, in order to generate correct predictions, the agent must grasp the interplay between these factors.

#### 2.2.2 Disentanglement Assuming Causal Interactions Among the Latent Factors

This section discusses the works that assume the latent factors (i.e. the causal and non-causal factors) have a causal interaction amongst each other. Majority of the works assume the non-causal features to act as confounding factors and aim to leverage front-door or back-door criterion to mitigate confounding bias and improve generalization. There also exist works that aim to learn the nature of the relationships among the latent factors. For instance, Yang et al. [31] suggest that when it comes to disentangling the generative elements of observational data, the bulk of frameworks presume that the latent factors are independent of one another. However, in practice, this assumption may not hold true. To solve this issue, they offer CausalVAE, a new framework. CausalVAE employs a Structural Causal Model layer that enables the model to recover latent components with semantics and structure using a Directed Acyclic Graph (DAG). First, the input x is encoded to get the independent components z. The collected factors are subsequently processed by the Structural Causal Model (SCM) layer, which converts independent factors into causal endogenous ones. Finally, a masking method is employed to replicate the assignment operation of SCMs by propagating the influence of parental variables on their children. Then the representation of the latent exogenous independent variables can be written as,

$$\epsilon = h(x, u) + \zeta \tag{18}$$

where  $\zeta$  is the independent noise term. After obtaining  $\epsilon$ , the causal representation is obtained as a linear SCM,

$$\mathbf{z} = \mathbf{A}^T \mathbf{z} + \epsilon = (I - \mathbf{A}^T)^{-1} \epsilon \tag{19}$$

where **A** represents the adjacency matrix of the causal graph, where  $\mathbf{A}_i \in \mathbb{R}^n$  is the weight vector such that  $\mathbf{A}_{ji}$  encodes the causal strength from  $z_j$  and  $z_i$ , and I is an identity matrix. Finally, the causal representations undergo a masking mechanism. This step resembles an SCM, which depicts how children are generated by their corresponding parental variables. The masking mechanism is formulated as,

$$z_i = g_i(\mathbf{A}_i \circ \mathbf{Z}; \boldsymbol{\eta}_i) + \epsilon_i \tag{20}$$

where  $\circ$  represents the element-wise multiplication,  $\eta_i$  is a parameter of the nonlinear function  $g_i$ . The final loss function consists of multiple terms, including the ELBO loss for VAE [32] and different constraints to ensure the latent variables **A** and z are identifiable.

#### Mitigating Confounding Bias via Backdoor Adjustments

The backdoor adjustment aids in approximating the interventional distribution, which can guide in identifying the causal link between the causal features and the image label. Let CS denote the representation of the causal features, S denote the representation of the non-causal features, and S denote the image label. The causal-graph under this scenario can be visualized as shown in Fig.5a. Similar to Eq.1, the backdoor adjustment can be formulated as,

$$P(Y \mid do(S = s_k)) = \sum_{i=1}^{|CS|} P(Y \mid S = s_k, CS = c_i) P(CS = c_i)$$
(21)

Zhang et al. [33] aim to utilize the causal graph to shed light on the poor generalization of classic person re-identification models when applied to unknown contexts. The primary premise of this research is that person pictures are influenced by two sets of latent random variables, namely identity-specific and domain-specific aspects. The authors propose Multi-Domain Disentangled Adversarial Neural Networks (MDANN), which learn two encoders from various datasets for embedding identity-specific (causal) and domain-specific (non-causal) components. To eliminate domain (identity) relevant information from embedded identity (domain) specific representations, the adversarial learning principle is used. Then, as shown in Eq. 21, a backdoor adjustment block (BA) is presented, which uses the identity-specific and domain-specific representations to achieve the approximation. The objective function is a combination of the backdoor adjustment and the classification loss.

Simiarly, Deng et al. [34] Attempt to use causality to enhance knowledge distillation among teacher-student models in order to improve generalizability. The authors create a causal graph to represent the causal linkages between the pre-trained instructor, the samples, and the prediction. To eliminate biased knowledge based on backdoor adjustment, the model employs the softened logits learnt by the teacher as the context information of an image. Overall, the proposed methodology captures full teacher representations while removing bias through causal intervention. Wang et al. [35] present a causal attention model capable of distinguishing between causal and confounding picture aspects. The authors use backdoor adjustment to accomplish disentanglement of the causative and confounding aspects. To address the over-adjustment problem, they create data splits repeatedly and gradually self-annotate the confounders.

### Mitigating Confounding Bias via Front-Door Adjustments

When the spurious features or confounders are unidentifiable or their distributions are hard to model, then one can use the front-door adjustment as it does not require explicitly modeling for the confounders. By introducing an intermediate variable Z between the input X and output Y, the front-door criterion transfers the requirement of modeling the intervening effects of confounders to modeling the intervening effects of the input. The motivation for this is clear if we consider a two state intervention. If we set the value of X, we can determine the corresponding value of Z, and we can then intervene again to fix that value of Z. By doing this for every value of Z we are able to determine the effect of X on Y. The causal-graph under this scenario can be visualized as shown in Fig.5b. The front-door criterion can be formulated as shown in Eq.2, where x' denotes the set of training data. By leveraging the front-door criterion, Li et al. [36] present an approach for mitigating confounding bias in the absence of identifying the confounders The proposed technique simulates the interventions among various samples using the front-door criteria and then optimizes the global-scope intervening impact on the instance-level interventions. Unlike previous studies, this is the first effort to use the front-door criteria for learning causal visual cues by taking the intervention among samples into account.

## 2.3 Invariance via Learning/Transfering Causal Mechanisms

These papers train mechanisms (such as neural networks) invariant across different domains. The conditionals in these cases remain invariant. These frameworks impose that the invariance of conditionals is valid as long as the conditionals represent the causal mechanisms. Formally, given  $\mathbf{X}$  as the input, Y as the label and  $S^*$  as the subset of invariant predictors, the conditional distribution  $P(Y^k \mid \mathbf{X}_{S^*}^k)$  is invariant across the k domains.

#### 2.3.1 IRM and its Extensions

When the training data originates from multiple domains, dividing the features into causal and non-causal features becomes challenging. One of the pioneering works - Invariant Risk Minimization (IRM) [4], was the first to address this problem. The authors propose to find a data representation such that a classifier trained on top of that representation matches for all domains. If the training data  $D_d := \{(x_i^d, y_i^d)\}_{i=1}^{n_d}$  is collected under multiple domains  $d \in \mathcal{E}_{\mathrm{tr}}$ . Then, to achieve high generalizability, the model should minimize the following loss,

$$R^{\text{OOD}}(f) = \max_{e \in \mathcal{E}_{\text{all}}} R^e(f)$$
 (22)

where  $R^e(f) := \mathbb{E}_{X^e,Y^e}[\ell(f(X^e),Y^e)]$  is the risk under domain e,  $\mathcal{E}_{all}$  is a large set of unseen domains that are related to the source domains,  $\ell(f(X^e))$  is the representation function for the input X from domain e, and  $Y^e$  is the true image label. IRM is a learning paradigm to estimate data representations eliciting invariant predictors  $w \circ \Phi$  across multiple domains. The constrained optimization problem that IRM aims to solve can be formulated as,

$$\min_{\boldsymbol{w}, \Phi} \sum_{e \in \mathcal{E}_{train}} R^{e}(\boldsymbol{w} \circ \Phi)$$
s.t.  $\boldsymbol{w} \in \underset{\hat{\boldsymbol{w}}}{\operatorname{argmin}} R^{e}(\hat{\boldsymbol{w}} \circ \Phi),$ 
(23)

where  $R^e$  is the cross-entropy loss for domain e,  $\Phi$  is the feature extractor and w is a linear classifier. Since this is a bi-level optimization problem, it is difficult to optimize. By adopting the first-order approximation, the loss function is,

$$\min_{\Phi} \sum_{e \in \mathcal{E}_{train}} R^e(\Phi) + \lambda \cdot \|\nabla_{w|w=1.0} R^e(w \circ \Phi)\|, \tag{24}$$

where  $w \in \mathbb{R}$  is a dummy classifier. To find the invariant features across different domains, the authors assume that the data from all the environments share the same underlying Structural Equation Model (SEM).

Ahuja et al. [37] show that IRM fails in settings where invariant features capture all information about the label in the input. They further show that along with the information bottleneck constraints, the invariance principle works in both settings – when invariant features capture the label's information entirely and when they do not. Krueger et al. [38] argue that the IRM principle is less effective when different domains are noisy. Furthermore, the authors also claim that IRM fails to achieve robustness w.r.t. covariate shifts. Also, the authors show that reducing differences in risk across training domains can reduce a model's sensitivity to a wide range of extreme distributional shifts, including the challenging setting where the input contains causal and anti-causal elements. Since the proposed model seeks robustness to whichever forms of distributional shift, it is more focused on the domain generalization problem than IRM. Finally, the authors prove that equality of risks can be a sufficient criterion for discovering causal structure. Li et al. [39] argue that the IRM principle fails for nonlinear classifiers and when pseudo-invariant features and geometric skews exist. To solve the problems of IRM, the authors in this work aim to utilize mutual information for causal prediction. Furthermore, they aim to adopt the variational formulation of the mutual information for nonlinear classifiers to develop a tractable loss function. The proposed model seeks to minimize invariant risks while at the same time mitigating the impact of pseudo-invariant features and geometric skews. Guo et al. [40] argue that the IRM principle guarantees the

existence of an invariant optimal classifier for a set of overlapping feature representations across domains. However, as DNNs tend to learn shortcuts, they can circumvent IRM by learning non-overlapping representations for different domains. Thus, IRM fails when the spurious correlations are stronger than the invariant relations. The authors propose to utilize conditional distribution matching to overcome this problem.

Huh et al. [41] propose that by conserving the class-conditioned feature expectation  $\mathbb{E}_e[f(x) \mid y]$  across the different domains, one could overcome the flaws in IRM. Bellot et al. [42] argue that earlier works that optimize for worst-case error under interventions on observed variables are typically not optimal under moderate interventions, especially in the presence of unobserved confounders. They also argue that minimum average error solutions can optimize for any dependency in the data, and their performance deteriorates when the data is subjected to interventions on observed variables; but better adjusts to interventions on the confounders. Thus, causal and minimum average error solutions can be interpreted as two extremes of a distributionally robust optimization problem with a range of intermediate solutions that may have a more desirable performance. To this end, the authors propose Derivative Invariant Risk Minimization (DIRM) that interpolates between the causal solution and the minimum average error solution.

The IRM algorithm has also been effectively utilized in various applications. For instance, Francis et al. [43] highlights the application of generalization in federated learning models. The authors argue that the current federated learning models have poor generalizability as most of the data for the federated server are not i.i.d making it difficult for these models to generalize better. To this end, they propose to leverage causal learning to improve the generalizability in a federated learning scenario. The proposed model consists of client and global server layers. The client layer extracts the features from their respective input data. The global layer aids the participating clients in exchanging intermediate training components and trains the federated model in collaboration by minimizing the empirical average loss. The authors leverage IRM to learn invariant predictors in a federated learning setup that can attain an optimal empirical risk on all the participating client domains.

### 2.3.2 Utilizing Auxiliary Functions to Model Conditional Distributions

Generally, the domain generalization problem can be solved by posing it from a causal discovery point of view. Muller et al. [44] argue that when posed from a causal discovery viewpoint, a set of features exists for which the relation between this set and the label is invariant across all domains. They claim that this happens because of the Independent Causal Mechanisms (ICM) principle, which states that every mechanism acts independently of the others. Given a joint distribution of the input **X**, the chain rule can decompose this distribution into a product of conditionals. This can be formulated as.

$$p_{\mathbf{X}}(x_1, \dots, x_D) = \prod_{i=1}^{D} p_i(x_i \mid \mathbf{x}_{pa(i)}),$$
 (25)

where  $\mathbf{x}_{pa(i)}$  refer to the causal parents of  $x_i$ , and the conditionals  $p_i$  of this causal factorization are called *causal mechanisms*. The normalizing flows model complex distributions using invertible functions T, which map the densities of interest to latent normal distributions. Thus, the authors represent the conditional distribution  $P(Y \mid h(\mathbf{X}))$  using a conditional normalizing flow. They seek to learn a mapping  $R = T(Y; h(\mathbf{X}))$  such that,  $R \sim \mathcal{N}(0, 1) \perp h(\mathbf{X})$  when  $Y \sim P(Y \mid h(\mathbf{X}))$ . T can be learned by minimizing the negative log-likelihood as,

$$\mathcal{L}_{\text{NLL}}(T, h) := \mathbb{E}_{h(\mathbf{X}), Y}[\|T(Y; h(\mathbf{X})\|^2 / 2 \\ \cdot -\log|\det \nabla_y T(Y; h(\mathbf{X}))|] + C,$$
(26)

where C is a constant,  $\det \nabla_y$  is the Jacobian determinant. Finally, the authors propose a differentiable two-part objective function formulated by,

$$\arg\min_{\theta,\phi} (\max_{e \in \mathcal{E}} \{ \mathcal{L}_{\text{NLL}}(T_{\theta}, h_{\phi}) \} + \lambda_I \mathcal{L}_I(P_R, P_{h_{\phi}(\mathbf{X}), E})), \tag{27}$$

where  $\phi$  and  $\theta$  denote model parameters,  $h_{\phi}$  denotes the feature extractor, E denotes the domain, and  $\mathcal{L}_{I}$  denotes the Hilbert Schmidt Independence Criterion (HSIC) [45], a kernel-based independence measure that penalizes dependence between the distributions of R and  $(h_{\phi}(\mathbf{X}), E)$ . The distribution loss stems from the theorem that if  $R \perp (h(\mathbf{X}), E)$ . Then, it holds that  $Y \perp E \mid h(\mathbf{X})$ .

### 2.3.3 Graphical Criterion based Methods

Zheng et al. [46] leverage the causal factorization from the ICM principle to identify the autonomous generating factors (i.e., distribution of each variable given its cause), as changing one will not affect others. Due to this autonomy, it is necessary to have varied factors to account for the distributional shifts. The authors classify this set of variables as the mutable set, which is not assumed to be known. The authors further argue that obtaining a set of invariant predictors is

possible, conditioned on the intervened mutable variables and any stable subset. However, it becomes a challenging task when the degeneration condition does not hold. The degeneration condition implies that the predictor with the whole stable set is min-max optimal if the intervened distribution can degenerate to the conditional one. When this condition is determined not to hold, the authors transform the worst-case quadratic loss into an optimization problem over all subsets of the stable set, the minimizer of which is sufficient and necessary for the predictor to be min-max optimal.

Lv et al. [47] propose a novel method, Causality Inspired Representation Learning (CIRL), which aims to learn causal representations. They argue that the intrinsic causal mechanisms (formalized as conditional distributions) can be feasible to construct if the causal factors are given. However, it is hard to recover the causal factors since they are unobservable. To alleviate this problem, they propose to learn representations based on three general properties of the causal factors. They are (1) The causal factors are separated from the non-causal factors; (2) The factorization of the causal factors should be jointly independent; and (3) the factors should be causally sufficient, i.e., they should contain the necessary causal information needed for the classification task. These properties are based on Common Cause Principle and Independent Causal Mechanism principle. First, CIRL performs interventions using Fourier transforms to differentiate between causal and non-causal features. Next, representations of the augmented and the original images are learned as  $r = \hat{g}(x)$ , where g(.) is the representation function. To simulate the causal factors that remain invariant to the intervention, the model enforces  $\hat{g}$  to keep unchanged dimension-wisely, as,

$$\max_{\hat{g}} \frac{1}{N} \sum_{i=1}^{N} COR(\tilde{\boldsymbol{r}}_{i}^{o}, \tilde{\boldsymbol{r}}_{i}^{a}), \tag{28}$$

where  $\tilde{r}_i^o$  and  $\tilde{r}_i^a$  are the normalized Z-score of the  $i^{th}$  column of  $\mathbf{R}^o = [(\mathbf{r}_1^o)^T, \dots, (\mathbf{r}_B^o)^T]^T \in \mathbb{R}^{B \times N}$  and  $\mathbf{R}^a = [(\mathbf{r}_1^a)^T, \dots, (\mathbf{r}_B^o)^T]^T$ , respectively, COR is a correlation function, and N is the number of dimensions. The joint independence among the causal factors is imposed by ensuring that any two dimensions of the representations are independent. Thus, the same dimension of  $\mathbf{R}^o$  and  $\mathbf{R}^a$  are considered positive pairs. In contrast, the different dimensions are considered negative pairs. Positive pairs should hold a stronger correlation when compared to negative pairs. This can be formulated with a factorization loss  $\mathcal{L}_{Fac}$  which can be formulated as,

$$\mathcal{L}_{Fac} = \frac{1}{2} \| \boldsymbol{C} - \boldsymbol{I} \|_F^2, \tag{29}$$

where C is the correlation matrix computed as,

 $C_{ij} = \frac{\langle ilde{r}_i^c, ilde{r}_j^a \rangle}{\| ilde{r}_i^c \| \| ilde{r}_j^a \|}, i,j \in 1,2,\ldots,N,$  and I is the identity matrix. Classification loss with the causal representations can be computed to ensure the representations are causally sufficient. However, inferior dimensions may carry relatively less causal information and make a small contribution to classification. Thus, the authors design an adversarial mask module to detect the inferior dimensions. The module learns the importance of the different dimensions. The dimensions corresponding to the largest  $\kappa \in (0,1)$  ratio are regarded as superior. In contrast, the rest are regarded as inferior ones. The authors use the GumbelSoftmax trick [16] to sample a mask with  $\kappa N$  values approaching 1. Thus, the classification loss is given by,

$$\mathcal{L}_{cls}^{sup} = \ell(\hat{h}_1(r^o \odot m^o), y) + \ell(\hat{h}_1(r^a \odot m^a), y) 
\mathcal{L}_{cls}^{inf} = \ell(\hat{h}_2(r^o \odot (1 - m^o)), y) + \ell(\hat{h}_2(r^a \odot (1 - m^a)), y).$$
(30)

The final objective function is represented as,

$$\min_{\hat{g},\hat{h}_1,\hat{h}_2} \mathcal{L}_{cls}^{sup} + \mathcal{L}_{cls}^{inf} + \tau \mathcal{L}_{Fac}, \quad \min_{\hat{w}} \mathcal{L}_{cls}^{sup} - \mathcal{L}_{cls}^{inf}, \tag{31}$$

where  $\tau$  is the trade-off parameter. Similarly, Wang et al. [48] argue that rather than the distribution of the features, it is the causal mechanism that remains invariant across domains. By utilizing feature-level invariance with the aid of regularizers, the model may learn spurious features that are only correlated to the label. Thus, to this end, the authors propose to treat the machine learning models as SCMs. In other words, the authors propose a novel constraint that measures the Average Causal Effect (ACE) between the causal attributions in the networks. They leverage contrastive learning and introduce ACE contrastive loss that pairs samples from the same class and different domains as positive pairs and those of different classes as negative pairs. The ACE contrastive loss regularizes the learning procedure and encourages domain-independent attributions of extracted features. The causal attribution of the neuron  $z^j$ , corresponding to the  $j^{th}$  feature, on the output y is calculated as the subtraction between the interventional expectation of y when  $z^j = \alpha$  and a baseline of  $z^j$ , given by,

$$c_{do(z^{j}=\alpha)}^{y} = \mathbb{E}[y \mid do(z^{j}=\alpha)] - \mathbb{E}_{z^{j}}[\mathbb{E}[y \mid do(z^{j}=\alpha)]]. \tag{32}$$

The ACE vector for the  $i^t h$  samples is given by,

$$\mathbf{c}_i = [c_i^1, c_i^2, \dots c_i^n].$$
 (33)

The final objective function is then formulated as,

$$\mathcal{L}_{\mathcal{A}}(\mathbf{D}; \theta, \phi) = \sum_{i=1}^{m} \ell(g_{\phi}(f_{\theta}(x_i)), y_i) + \rho \sum_{i=1}^{m} \mathcal{A}_{\theta, \phi}(x_i),$$
(34)

where the first term denotes the classification loss,  $g_{phi}$  is the classifier and  $f_{\theta}$  are the features and  $\mathcal{A}_{\theta,\phi}(x_i)$  represents the Contrastive ACE loss which is given by,

$$\max\{\operatorname{dist}(\mathbf{c}_i, \mathbf{c}_{q(\mathcal{P}_i)}) - \operatorname{dist}(\mathbf{c}_i, \mathbf{c}_{q(\mathcal{N}_i)}) + \delta, 0\}, \tag{35}$$

where  $q(\mathcal{P}_i)$  is an index sampled uniformly at random from the positive pairs' set,  $q(\mathcal{N}_i)$  is from the negative pair,  $\delta$  is a margin variable, dist is a distance metric. Sun et al. [49] argue that when dealing with sensory-level data such as modeling pixels, it is beneficial to model the problem similar to human perception; i.e., the causal factors of the label Y is related to unobserved abstractions S via a mechanism  $f_y$  such that  $Y \leftarrow f_y(S, \varepsilon_y)$ , where  $\varepsilon$  is a noise term. At the same time, there exist latent variables Z that, along with variables S generate the input image X via mechanism  $f_x$  such that  $X \leftarrow f_x(S, Z, \varepsilon_x)$ . Under this situation, domain shifts occur when variables Z are allowed to correlate to the variables S spuriously. For instance, when dealing with the image classification problem, the background features can be classified as S, and the object-related abstractions such as shape can be classified as S. The authors encapsulate this information in a set of causal models. They argue that the generating mechanisms  $f_x$  and  $f_y$  are invariant across domains. At the same time, the spurious relation between S and S is allowed to vary. Mathematically, Causal Invariance refers to the condition when S is an experience of the variational Bayesian method to estimate the Causal Invariance during training and optimize it during testing.

### 2.3.4 Kernel-based Optimization Methods

Muandet et al. [50] was one of the earliest works aiming at enhancing the generalizability of the machine learning models from a causal perspective. The authors argue that the conditional distribution of the label Y, given an input X is stable. However, the marginal distribution i.e., P(X) may fluctuate smoothly. Due to this fluctuation machine learning models may suffer from model misspecification, i.e. the model fails to account for everything that it should. To alleviate this problem, the authors propose Domain Invariant Component Analysis (DICA). DICA aims to find data transformations that minimize the difference between the marginal distribution of different domains while preserving the stable conditional  $P(Y \mid X)$ . They introduce a domain generalization approach that learns an invariant transformation across domains between inputs and outputs by minimizing the dissimilarities between different domains. Basically, the objective of this work is to find a transformation that satisfies the following two properties: (1) Minimizing the distance between the distribution of the samples transformed via this transformation; and (2) The learned transformation between input and output remains invariant across different domains. To do so, a kernel-based optimization objective is defined as:

$$\max_{B \in \mathbb{R}^N \times M} \frac{\frac{1}{n} \operatorname{Tr}(B^T L (L + n\epsilon I_n)^{-1} K^2 B)}{\operatorname{Tr}(B^T K Q K B + B K B)},$$
(36)

where K and Q are the block kernel and coefficient matrices, respectively and B is the estimator which satisfies the two desired properties.

### 3 Domain Generalization in Text and Graphs

So far we have surveyed existing works on the causal domain generalization for the image data. In the following we discuss the frameworks designed for text and graph data.

# 3.1 Invariance Learning for Text

With the advances of large pre-trained models, Natural Language Processing models have gained widespread success over multiple applications in the real world. However, these models are brittle to out-of-domain samples [51]. A series of works, showcase how the language models rely on spurious correlations for classification. For instance, Wang et al. [52] show that words such as *Spielberg* is correlated to positive movie reviews. Wulczyn et al. [53] show that a toxicity classifier learns that "gay" is correlated with toxic comments. To address the reliance on spurious correlations recent works [54, 55, 56, 57, 58] aim at leveraging different causal techniques to improve the generalizability of the NLP models. Interested readers can refer to [59] which focuses on how to infer causal relations in natural language processing models.

Causality-aware domain generalization works in the natural language domain are also aligned with our defined categories. For instance, recent works in NLP [60, 61] leveraged human in-the-loop systems to make the models robust to spurious correlations by leveraging human common sense of causality. They augment training data with crowd-sourced annotations about reasoning of possible shifts in unmeasured variables and finally conduct robust optimization to control worst-case loss. While these works highlight that human annotations improve the robustness, collecting such annotations can be costly. To overcome this problem Wang et al. [55] propose to train a robust classifier with automatically generated counterfactual samples which is aligned with causal data augmentation approaches as listed earlier in the categorization section. The authors aim at utilizing the closest opposite matching approach to identify the likely causal features, and then generate counterfactual training samples by substituting causal features with their antonyms and assigning opposite labels to the newly generated samples.

Mathematically, for a given sample (d, y), where d represents document and y represents the label, the corresponding counterfactual (d', y') is generated by (i) substituting causal terms in d with their antonyms to get d', and (ii) assigning an opposite label y' to d'. The authors propose a two-stage process as follows:

- Using models such as logistic regression, the authors identify the strongly correlated terms  $\langle t_1 \dots t_k \rangle$  as candidate causal features.
- For each top term t and a set of documents containing  $t: D_t = \langle d_1 \dots d_n \rangle$ , the authors utilize context similarity to identify documents with *opposite* labels.
- Then, the authors select the highest score match for each term and identify likely causal features by picking those whose closest opposite matches have scores greater than a threshold (0.95 is used below).
- Then for each sample, the authors substitute the causal terms with antonyms and generate counterfactual samples. Finally, they train the robust classifier using the original and counterfactual data.

Similarly, Tan et al. [62] propose to construct meaningful counterfactuals that would reflect the model's decision boundaries. The authors leverage rule-based schemes that negate causal relations or strengthen conditionally causal sentences.

Choi et al. [54] aim to utilize contrastive learning to enhance the representations of the causal features, this is aligned with capturing invariance via learning causal representations. The porposed model  $C^2L$ , first aims to identify the causal tokens based on attribution scores. Formally, to identify the important tokens, the authors leverage attribution scores as follows:

$$g_i = \|\nabla_{\mathbf{w}_i^p} \mathcal{L}_{\text{task}}(x, y; \phi)\|^2, \tag{37}$$

where x denotes the input, y denotes the label,  $g_i$  denotes the gradient magnitude computed from the classifier  $f_{\phi}$ , and  $\mathcal{L}_{task}$  denotes the cross-entropy loss. The gradient-based score of token w is aggregated over all the training texts having the token w by:

$$s^{grad}(w) = \sum_{(x,y)\in\mathcal{D}} \frac{1}{n_{w,x}} \sum_{i\in\{1,\dots,T\}} \mathbb{I}(w_i = w) \cdot g_i$$
 (38)

where  $\mathbb{I}$  is an indicator function, and  $n_{w,x}$  is the number of word w in the input x. After obtaining the scores for each tokens, the authors employ a causal validation technique to identify the causal tokens. As per this step, the authors compute the Individual Treatment Effect (ITE) of each token on the label. The main intuition behind this step is that, if the masked text can be reconstructed into multiple examples with different classes, we can decide the masked term has a causal effect. To this end, the authors use BERT with dropout mechanism to identify the top-k substitutions for the token w. The k candidates are then passed through the classifier to obtained the predicted labels  $\hat{y}$ . By testing whether the k labels are evenly distributed into the classes, we can decide the high-attributed token w as causal to its task label y.

Finally, the authors leverage contrastive learning to better learn the causal structure of the classification task. After obtaining the causal features, the authors generate causal triplets of the form  $(x, x^+, x^-)$ .  $x^-$  denotes the counterfactual pair that is generated by masking out causal words. In contrast,  $x^+$  denotes the factual pair generated by masking one of the non-causal words that is still recognized as the original label y, which helps to learn a model invariant to these features. The contrastive objective aims at mapping the representation of x closer to  $x^+$  and further from  $x^-$ . The objective can be formulated as,

$$\mathcal{L}_{c}(x;\theta) = \max(0, \Delta_{m} + \frac{1}{J} \sum_{j=1}^{J} s_{\theta}(x, x_{j}^{+}) - \frac{1}{J} \sum_{j=1}^{J} s_{\theta}(x, x_{j}^{-}))$$
(39)

where J is the number of positive/negative pairs,  $\Delta_m$  is a margin value and  $s_{\theta}(\cdot, \cdot)$  is distance between the representations. The final objective function is given by,

$$\mathcal{L} = \mathcal{L}_{\text{task}} + \lambda \mathcal{L}_c \tag{40}$$

where  $\lambda$  is a balancing coefficient for the contrastive objective.

Inspired by IRM [4], Peyrard et al. [57] propose invariant Language Modeling (iLM), a framework that generalizes better across multiple environments for language models in NLP. iLM is aligned with learning invariance by transferring causal mechanisms. The ILM aims at utilizing a game-theoretic implementation of IRM for language models. In this case, the invariance is achieved by a specific training schedule in which each environment competes with the others to optimize their environment-specific loss by updating subsets of the Language Model (LM) in a round-robin fashion.

The IRM-games [63] method aims to change the training procedure of IRM by using a game-theoretic perspective in which each environment e is tied to its own classifier  $w^e$ , and the feature representation  $\phi$  is shared. The global classifier w is defined as an ensemble formulated as,

$$w = \frac{1}{|\mathscr{E}|} \sum_{e \in \mathscr{E}} w^e \tag{41}$$

where  $\mathscr E$  represents the set of training environments. Each environment takes turns to minimize their own empirical risk  $R^e(w \circ \phi)$  w.r.t their own classifier  $w^e$ , while the shared  $\phi$  is updated periodically. The authors adopt this IRM-games setup for language models. First, the shared representation can be represented by the main body of the encoder of an LM, and  $w^e$  is the language modeling head that outputs the logits after the last layer. The multiple environments for this task can be the different sources from which text data emerges, for instance, encyclopedic texts, Twitter, news articles, and so on. Suppose for each environment the data can be represented as  $\{(X^e,Y^e)\}_{e=1...n}$ . A forward pass on a batch  $(x_i,y_i)$  sampled according to  $P(X^i,Y^i)$  from environment i involves n language modeling heads  $\{w_e\}_{e=1...n}$ :

$$\hat{y} = \operatorname{softmax}(\sum_{e=1}^{n} w_e \circ \phi(x_i))$$
(42)

Based on the task, a modeling loss  $\mathcal{L}$  can be applied to the output  $\hat{y}$ . Since the underlying causal model is not known for language models, utilizing the data stemming from different environments can facilitate in learning the invariant relationships. However, the choice of environments is an important step as the environments define which relations are spurious in nature.

#### 3.2 Invariance Learning for Graphs

Unlike vision and natural language data, graph data is heterogeneous. Graph Neural Networks (GNNs) fuse heterogeneous information from node features and graph structures to learn effective node embeddings. However, the complex and unobserved non-linear dependencies among representations are much more difficult to be measured and eliminated than the linear cases for decorrelation of non-graph data. In out-of-distribution scenarios, when complex heterogeneous distribution shifts exist, the performance of current GNN models can degrade substantially, mainly induced by spurious correlations. The spurious correlations intrinsically come from the subtle correlations between irrelevant and relevant representations. To address this problem, a variety of methods have been proposed.

Liu et al. [64] leverage graph data augmentations to identify graph rationales. The authors aim to augment the rationale subgraph by removing its environment subgraph and combining it with different environment subgraphs. Furthermore, they propose a framework, namely, GREA, that leverages masking to separate the rationales from the environment. After learning the node representations with the aid of a GNN, a Multi-Layer Perceptron (MLP) is used to map the node representations to a mask vector  $\mathbf{m} \in (0,1)^N$  on the nodes in the set  $\mathcal{V}$ .  $m_v = \Pr\left(v \in \mathcal{V}^{(r)}\right)$  is the node-level mask that indicates the probability of node  $v \in \mathcal{V}$  being classified into the rationale subgraph. It is formulated as,

$$m = \sigma \left( \text{MLP}_1 \left( \text{GNN}_1(g) \right) \right) \tag{43}$$

GREA uses another GNN encoder to generate contextualized node representations  $\mathbf{H} : \mathbf{H} = \text{GNN}_2(g)$ . With  $\mathbf{m}$  and  $\mathbf{H}$ , the rationale subgraph and environment subgraph can be easily separated in the latent space. The rationale and environment subgraph are generated as,

$$\mathbf{h}^{(r)} = \mathbf{1}_N^\top \cdot (\mathbf{m} \times \mathbf{H}), \quad \mathbf{h}^{(e)} = \mathbf{1}_N^\top \cdot \left( (\mathbf{1}_N - \mathbf{m}) \times \mathbf{H} \right),$$

where  $\mathbf{1}_N$  denotes the N-size column vector with all entries as 1, and  $\mathbf{h}^{(r)}, \mathbf{h}^{(e)} \in \mathbb{R}^d$  are the representation vectors of rationale graph  $g^{(r)}$  and environment graph  $g^{(e)}$ , respectively. After the rationale and environment separation, the model leverages two augmentation strategies to make predictions. First, it combines each rationale subgraph with multiple environment subgraphs to generate augmented samples which can improve the model's robustness and generalization. The prediction is made as,

$$\hat{y}_{(i,j)} = MLP_2\left(\mathbf{h}_{(i,j)}\right) \tag{44}$$

where  $\hat{y}_{i,j}$  is computed as,

$$\mathbf{h}_{(i,j)} = AGG\left(\mathbf{h}_i^{(r)}, \mathbf{h}_j^{(e)}\right) = \mathbf{h}_i^{(r)} + \mathbf{h}_j^{(e)}$$

$$\tag{45}$$

Second, it removes the environment subgraph and uses only the rationale subgraphs to make predictions as follows,

$$\hat{y}_i^{(r)} = \text{MLP}_2\left(\mathbf{h}_i^{(r)}\right) \tag{46}$$

For instance, Fan et al. [65] propose to leverage causality to overcome the subgraph-level spurious correlations to improve the GNN generalizability. They analyzed the degeneration of GNNs from a causal view and propose a novel causal variable distinguishing regularizer to decorrelate each high-level variable pair by learning a set of sample weights. The sample reweighting method aids in eliminating the dependence between high-level variables, where non-linear dependence is measured by weighted Hilbert-Schmidt Independence Criterion (HSIC) [45]. HSIC is formulated as,

$$HSIC_0^{k,l}(U,V,\mathbf{w}) = (m-1)^{-2}\operatorname{tr}(\hat{\mathbf{K}}\mathbf{P}\hat{\mathbf{L}}\mathbf{P})$$
(47)

where  $\hat{\mathbf{K}}, \hat{\mathbf{L}} \in \mathbb{R}^{m \times m}$  are weighted RBF kernel matrices containing entries  $\hat{\mathbf{K}}_{ij} = k \left( \hat{U}_i, \hat{U}_j \right)$  and  $\hat{\mathbf{L}}_{ij} = l \left( \hat{V}_i, \hat{V}_j \right)$ ,  $\hat{U} = (\mathbf{w} \cdot \mathbf{1}^T) \odot U$ , and  $\hat{V} = (\mathbf{w} \cdot \mathbf{1}^T) \odot V$ . Finally, the weights are optimized as,

$$\mathbf{w}^* = \underset{\mathbf{w} \in \Delta_m}{\operatorname{arg \, min}} \sum_{1$$

where  $\Delta_m = \{\mathbf{w} \in \mathbb{R}^n_+ \mid \sum_{i=1}^m \mathbf{w}_i = m\}$ , and we utilize  $\mathbf{H}_{,0:d}$  denotes the treatment variable, and  $\mathbf{H}_{,(p-1)d:pd}$  denotes the confounders,  $\mathbf{w} = \operatorname{softmax}(\mathbf{w})$  to satisfy this constrain Hence, reweighting training samples with the optimal  $\mathbf{w}^*$  can mitigate the dependence between high-level treatment variable with confounders to the greatest extent. Sui et al. [66] propose a causal attention model that could distinguish between causal and confounding features of a graph. The authors leverage the backdoor adjustment to disentangle the causal and confounding features. The proposed model separates the causal and confounding features from full graphs using attention mechanism.

Bevilacqua et al. [67] leverage causal models inspired by Stochastic Block Models (SBM) [68] and graphon random graph models [69] to learn size-invariant representations that better extrapolate between test and train graph data. The authors construct graph representations from subgraph densities and use attribute symmetry regularization to mitigate the shift of graph size and vertex attribute distributions.

Chen et al. [70] propose to provide guaranteed OOD generalization on graphs under different distribution shifts. The authors leverage three SCMs to characterize the distribution shifts that could happen on graphs. They further argue that GNNs are invariant to distribution shifts if they focus only on a invariant and critical subgraph  $G_c$  that contains the most of the information in G about the underlying causes of the label. To learn the invariant representation, the authors propose to align with two causal mechanisms that occur during graph generation, i.e.,  $C \to G$  and  $(G_s, E_G, G_c) \to G$  where C represents the causal features, S represents the non-causal features,  $E_G$  denotes the environment,  $G_c$  inherits the invariant information of C that would not be affected by the interventions, and  $G_s$  inherits the varying features. The alignment is realized by decomposing a GNN into a featurizer GNN  $g: \mathcal{G} \to \mathcal{G}_c$  aiming to identify the desired  $G_c$ ; b) a classifier GNN  $f_c: \mathcal{G}_c \to \mathcal{Y}$  that predicts the label Y based on the estimated  $G_c$ , where  $\mathcal{G}_c$  refers to the space of subgraphs of G. Formally, the learning objectives of  $f_c$  and g can be formulated as:

$$\min_{f_c,g} R\left(f_c\left(\hat{G}_c\right)\right), \text{ s.t. } \hat{G}_c \perp E, \hat{G}_c = g(G),$$

where  $R\left(f_c\left(\hat{G}_c\right)\right)$  is the empirical risk of  $f_c$  that takes  $\hat{G}_c$  as innuts to predict label Y for G, and  $\hat{G}_c$  is the intermediate subgraph containing information about C and is independent of E. The final objective is given by,

$$\max_{f_c,g} I\left(\hat{G}_c; Y\right) + I\left(\hat{G}_s; Y\right),$$
s.t.  $\hat{G}_c \in \underset{G_c = g(G), \tilde{G}_c = g(\tilde{G})}{\arg \max} I\left(\hat{G}_c; \tilde{G}_c \mid Y\right),$ 

$$I\left(\hat{G}_s; Y\right) \leq I\left(\hat{G}_c; Y\right), \hat{G}_s = G - g(G),$$
(49)

This work differs from [67] as it establishes generic SCMs that are compatible with several graph generation models, and different types of distribution shifts.

Wu et al. [71] aims at identifying graph rationales that capture the invariant causal patterns. To this end, the authors propose Discovering Invariant Rationales (DIR), which leverages causal interventions to instantiate environments and further distinguish the causal and non-causal parts. The task of invariant rationalization can be formulated as,

$$\min_{h_{\tilde{C}},h_{\hat{Y}}} \mathcal{R}\left(h_{\hat{Y}} \circ h_{\tilde{C}}(G), Y\right), \quad \text{s.t. } Y \perp \tilde{S} \mid \tilde{C},$$

$$(50)$$

where  $h_{\tilde{C}}$  discovers rationale  $\tilde{C}$  from the observed G,  $h_{\hat{Y}}$  represents the classifier,  $\tilde{C}$  is the causal rationale and  $\tilde{S} = G \backslash \tilde{C}$  is the complement of  $\tilde{C}$ . Furthermore, to obtain the environments the authors generate s-interventional distribution by doing intervention do(S=s) on S, which removes every link from the parents PA(S) to the variable S and fixes S to the specific value s. By stratifying different values  $\mathbb{S} = \{s\}$ , they obtain multiple s-interventional distributions. The DIR Risk is formulated as,

$$\min \mathcal{R}_{\text{DIR}} = \mathbb{E}_s[\mathcal{R}(h(G), Y \mid do(S=s))] + \lambda \operatorname{Var}_s(\{\mathcal{R}(h(G), Y \mid do(S=s))\})$$
(51)

where  $\mathcal{R}(h(G), Y \mid do(S = s))$  computes the risk under the s-interventional distribution,  $Var(\cdot)$  calculates the variance of risks over different s-interventional distributions and  $\lambda$  is a hyper-parameter to control the strength of invariant learning.

# 4 Benchmark, Evaluation, and Code Repositories

In this section, we provide a comprehensive review of benchmark datasets and evaluation metrics for all three types of data, i.e., image, text and graph, for the causal domain generalization task.

### 4.1 Benchmark Datasets

Causality-aware domain generalization has been studied across various applications, including but not limited to computer vision, natural language processing, and graphs. Tables 1 and 2 summarize the commonly used datasets based on the different applications. In this section, we briefly describe these datasets and the applications.

**Face Detection** can be decomposed into multiple tasks, such as face attributes detection, and human face synthesis. Some of the benchmark datasets are CelebA [72] which contains auxiliary attribute labels (such as GENDER, SMILE) to improve the generalization performance. A range of datasets, including Face-Forensics++ [73] and DeeperForensics-1.0 [74], were leveraged to facilitate generalization for the human face synthesis task.

**Emotion detection** has been studied extensively to aid applications such as mental health care and driver drowsiness detection. Benchmark datasets such as CK+ [75], MMI [76], and Oulu-CASIA [77] have been leveraged to test generalization capabilities for emotion detection. These datasets contain multiple face angles along with basic expression labels such as *anger*, *disgust*, and *fear*.

Handwritten digit recognition are a good example of distribution shifts as the different writing styles of people are as different domains and the shape of the digit remains invariant. Many methods have leveraged benchmark digits datasets to solve the generalization problem in digits recognition. Some majorly applied datasets are FashionMNIST [78] (which is a collection of grayscale fashion article images), ColoredMNIST+ [40], MNIST-M [79], SVHN [80], and SYN [79], MNIST [3] and its variants such as ColoredMNIST [4] (where the digits have different color distributions) and RotatedMNIST [81] (where the digits are rotated on different angles).

**Medical Imaging** refers to an umbrella term comprising multiple tasks. In the medical scenario, domain shifts are mainly caused by different acquisition processes. Thus, to improve the generalization performance, a range of works [13, 19, 49] leveraged real-world medical imaging datasets such as Alzheimer's Disease Neuroimaging Initiative (ADNI) [82]. Furthermore, the Chest X-rays dataset contains images from three sources: NIH [83], ChexPert [84] and RSNA [85]. The central task for this dataset was to classify the image to whether the patient has Pneumonia (1) or not (0).

**Body pose estimation** refers to the 3D pose estimation task. Recently the authors et al. [17] proposed to leverage causality to improve the cross-domain pose estimation problem. They utilize multiple benchmark datasets such as Human3.6M [86], 3DPW [87], MPI-INF-3DHP (3DHP) [88], SURREAL [89], and HumanEva [90] for body pose estimation.

**Person Re-IDentification** aims at matching person images of the same identity across multiple camera views. In this scenario, the domain shift arises in image resolution, viewpoint, lighting condition, background, etc. Some benchmark

| Dataset            | Description  | Domain | Downstream<br>Task        | Basic Statistics<br>of the Dataset   | Useful in which<br>Category and How?   |
|--------------------|--|--------|---------------------------|--|--|
| CelebA             | is a large-scale face attributes dataset.<br>The images in this dataset cover large<br>pose variations and background clutter.   | Vision | Face<br>Detection         | 202,599 celebrity images<br>10,177 unique celebrities<br>40 attribute annotations. | Since it provides access to auxiliary labels<br>such as GENDER and SMILE, it has<br>been utilized across multiple causal methods that<br>aim to learn and identify the causal features of an image.  |
| CK+                | The Extended Cohn-Kanade (CK+) dataset contains video sequences ranging from 18 to 50 years of age with various genders and heritage. The videos are labeled to denote the emotions such as anger, disgust, and so on.   | Vision | Emotion<br>Detection      | 593 video sequences<br>123 different subjects<br>327 labelled videos               | This dataset has been used for causal data augmentation by combining different features useful to predict emotions with confounding features such as noise.  |
| MNIST              | The MNIST database (Modified National Institute of<br>Standards and Technology database)<br>is a large collection of handwritten digits.   | Vision | Digit<br>Recognition      | 70,000 instances<br>10 classes   | MNIST is a part of the DigitsDG dataset that contains 4 domains, including MNIST, MNIST-M SVHN and SYH. It has been leveraged by causal methods trying to learn and identify the causal features.  |
| CMNIST             | Colored MNIST (CMNIST) is a synthetic<br>binary classification task derived from MNIST.<br>In CMNIST the color and the label have a different<br>correlation in the train set when compared to the<br>correlation in the test set.   | Vision | Digit<br>Recognition      | 70,000 instances<br>2 classes  | CMNIST was a synthetic dataset introduced in the<br>IRM paper. Since the spurrious correlation between<br>color and the label changes in the training and test<br>set, this dataset has been used in a wide-range of<br>works that learn invariance by transfering the<br>causal mechanisms. |
| Chest X-Ray14      | ChestX-ray14 is a medical imaging dataset<br>which contains frontal-view X-ray images<br>of unique patients along with disease<br>labels mined from text.  | Vision | Medical<br>Imaging        | 112,120 X-ray images<br>30,805 unique patients<br>23 years worth data.             | Chest X-ray14 is a part of the ChestX-rays dataset that contains 3 domains, including, ChestXray14, MIMIC-CXR, and Stanford CheXpert. These datasets have been leveraged in causal methods trying to learn and identify the causal features.   |
| Human3.6M          | Human3.6M is a 3D human pose and corresponding<br>images dataset. It contains various actors both<br>male and female posing in different scenarios.  | Vision | Human Pose<br>Estimation  | 3.6 million images<br>11 professional actors<br>17 scenarios                       | Since this dataset provides access to different poses under different domains, this dataset has been leveraged by causal data augmentation approaches where given the domain and the content, the models generate causally augmented images to train the model.                              |
| CUHK03             | The CUHK03 consists of images of different identities, where 6 campus cameras were deployed for image collection and each identity is captured by 2 campus cameras. This dataset provides two types of annotations, one by manually labelled bounding boxes and the other by bounding boxes produced by an automatic detector. | Vision | Person<br>Re-ID           | 14,097 images<br>1,467 identities  | Since this dataset reflects the same person across multiple backgrounds, this dataset has been leveraged by causal methods that aim to learn and identify the causal features present in an image.   |
| WaterBirds         | The WaterBirds dataset consists of water birds<br>and land birds. It was extracted from<br>Caltech-UCSD Birds-200-2011 benchmark dataset,<br>with water and land background<br>extracted from the Places dataset   | Vision | Birds<br>Classification   | ~5,000 images<br>2 classes   | Since this dataset contains auxiliary labels such<br>as the label for the background of an image,<br>they have been utilized in causal methods<br>that aim to learn and identify the causal<br>features present in an image.   |
| Terra<br>Incognita | This dataset contains images from twenty camera traps which were deployed to monitor animal populations. Since the traps are fixed, the background changes little across images. Capture is triggered automatically, thus eliminating human bias.  | Vision | Animals<br>Classification | 24,788 images<br>10 classes<br>4 domains   | Since this dataset contain high cross-domain discrepancies it has been used in causal methods that aim to learn and identify the causal features and by causal methods that capture invariance through the transfer of causal mechanisms   |
| PACS               | PACS refers to Photo, Art, Cartoon, and Sketch.<br>Each of the domain have seven categories.   | Vision | Object<br>Recognition     | 9,985 images<br>4 domains<br>7 categories  | Since this dataset contains the representation of the same object across multiple image styles, it allows models to distinguish between causal and non-causal features. Thus, it is widely used in causal methods that aim to learn and identify the causal features present in an image.    |
| OfficeHome         | Office-Home is a benchmark dataset that<br>represents images of the same object under<br>different scenarios, including Art, Clipart,<br>Product and Real-World.   | Vision | Object<br>Recognition     | 15,500 images<br>65 categories<br>4 domains  | Since this dataset contains the representation of the same object across multiple image styles, it allows models to distinguish between causal and non-causal features. Thus, it is widely used in causal methods that aim to learn and identify the causal features present in an image.    |

Table 1: A description of different vision benchmark datasets and how they have been leveraged by different causality-aware domain generalization methods.

| Dataset                  | Description   | Domain | Downstream<br>Task               | Basic Statistics<br>of the Dataset                                    | Useful in which<br>Category and How?  |
|--------------------------|---|--------|----------------------------------|---|---|
| PROTEINS                 | PROTEINS is a dataset of proteins that are classified as enzymes or non-enzymes.  Nodes represent the amino acids and two nodes are connected by an edge if they are less than 6 Angstroms apart. | Graph  | Graph<br>Classification          | 1,113 graphs<br>39 Avg. no. of nodes<br>2 classes                     | This dataset is used along with other biological graph datasets such as MUTAG and NCII by causal methods that aim to learn and identify the causal features present in an image.  |
| Multi-NLI                | The Multi-Genre Natural Language Inference<br>dataset is a collection of written and spoken<br>english data over ten different genres.  | Text   | Natural<br>Language<br>Inference | 433K sentence-pairs<br>10 domains                                     | This dataset contains spurious correlations as usually, the second sentence in a pair containing a negation is classified as contradiction. Works leveraging causal data augmentation have tried to tackle this spuriousness by generating counterfactuals for this situation.  |
| CivilComments            | Civil Comments contains the archive of the Civil Comments platform. It is annotated for toxicity across different demographic groups in the english language.                                     | Text   | Toxicity<br>Detection            | ~2 million comments<br>8 domains (demographic<br>groups)<br>2 classes | This dataset has been leveraged by causal methods that aim to learn and identify the causal features from text as it is possible to find the spurious correlations between different words and the toxicity labels of the comments.   |
| Amazon Kindle<br>Reviews | The Amazon Kindle Reviews dataset contains product reviews from kindle. The main goal is to infer the review's sentiment as positive or negative.   | Text   | Sentiment<br>Classification      | 10,500 reviews  | This dataset is leveraged by the causal data augmentation methods as these models perform data augmentations on non-causal (non-sentimental) words to study the domain generalization problem. It has also been used by causal methods that aim to learn and identify the causal features from text as the sentimental words have a causal relation to the sentiment classification task. |

Table 2: A description of different Graph and Text benchmark datasets and how they have been leveraged by different causality-aware domain generalization methods.

datasets for this task are CUHK02 [91], CUHK03 [92], Market1501 [93], DukeMTMC-ReID [94], CUHKSYSU PersonSearch [95].

**Animal and Birds Classification** refers to classifying different animals and bird species observed across multiple environments. For instance, the WaterBirds dataset contains images of water birds (Gulls) and land birds (Warblers) extracted from the Caltech-UCSD Birds-200- 2011 (CUB) dataset [96] with water and land background extracted from the Places dataset [97]. In addition, there are multiple datasets for animal classification, such as iWiLDSCam [98], and TerraIncognita [99].

**Object Recognition** is one of the most prominent tasks for studying the domain generalization problem, where the domain shift varies substantially across different datasets. For instance, PACS [100], OfficeHome [101], DomainNet [102] and ImageNet-R [103], deal with image style changes where the same object is varied across different styles. Another commonly used datasets are DomainNet [102], VLCS [104], ImageNet-C [105], NICO [106] and NICO++ [107].

**Graph Classification** is a crucial activity when dealing with graph distribution shifts. For this endeavor, a variety of datasets have been chosen to help with domain generalization. One such small-scale real-world dataset is HIV, which was modified from MoleculeNet [108]. ZINC is a real-world dataset for molecular property regression from the ZINC database [109]. Motif is a synthetic base-motif dataset motivated by Spurious-Motif [71].

**Node Classification** is another prominent task when dealing with node classification. Currently there exists several datasets that contain the out-of-distribution samples to test graph generalization capabilities, such as Cora [110], and Arxiv [111]. There also exist synthetic dataset such CBAS derived from the BA-Shapes [112].

**Sentiment Classification** deals with the task of classifying sentiments from human documents. Currently there exists several benchmark datasets for this task, examples include, Amazon Reviews Dataset [113], IMDb dataset [114], FineFood dataset [115]. Various works such as [54] aim to utilize multiple datasets listed here by training the model on one of the benchmark datasets and evaluate it against the other benchmarks.

**Toxicity Detection** refers to the task of detecting toxicity in textual data. One of the well-known benchmarks is the CivilComments dataset [116]. When evaluating for the cross-domain generalizability, the model is tested to see whether the model can detect toxicity without depending on the demographic identities.

**Natural Language Inference** aims at classifying pair of texts based on their logical relationship. One of the benchmark datasets for NLI is MultiNLI [117]. A series of works such as [54] have utilized this benchmark.

#### 4.2 Evaluation

Causality-aware domain generalization methods employ evaluation mechanisms that are very similar to standard domain generalization methods. Furthermore, the causal theories and methodologies employed by these frameworks help in the identification of invariant relationships while doing the same downstream task. Thus, the evaluation procedure is determined by the technique's nature, i.e., whether the approach is a single source domain generalization (where the model is trained on data from a single source) or a multi-source domain generalization (where the model is trained on data from multiple sources).

For instance, the *leave-one-domain-out rule* is one of the most prominent when dealing with multi-source domain generalization [13, 23, 31, 118]. Per this rule, given a dataset containing at least two distinct domains, multiple are used as source domains for training the model while one is used as the target domain, on which the model is directly tested without any adaptation. Another commonly employed strategy for domain generalization is *Training-domain validation set* [19, 119]. In this setting, the source domains are split into two parts; one is used for training while the other is used for validation. While training the model, each domain's training parts are combined, and the validation parts are used to select the best model. When dealing with single domain generalization methods, the popular approach is to train the model on the source domain and directly evaluate it on the different test domains.

Classification accuracy, top-1 error rate, top-5 error rate, and DICE scores are some of the main measures used to assess these models. The classification accuracy metric indicates how well the model classifies samples. The top-1 error rate determines if the model's projected top class corresponds to the target label. In the instance of the Top-5 mistake rate, we examine whether or not the target label appears in the top-5 predictions. Finally, the DICE scores compare the pixel-by-pixel agreement of a projected segmentation with its matching ground truth. Aside from these commonly used metrics, some works [20] leveraged mean Corruption Error (mCE) and mean relative Corruption Error (mrCE), which are commonly used to evaluate a model's robustness.

Although the causality-aware models strive to perform the same downstream job, we believe that causal evaluation procedures and metrics, in addition to the traditional setting, are required to verify the models' causal features. For example, Yang et.al [31] leveraged Maximal Information Coefficient (MIC) and Total Information Coefficient

(TIC) [120] as additional evaluation metrics. They use these metrics because they reflect the degree of information relevance between the learnt representation and the ground truth labels of ideas, which is important because their model tries to learn causal representations.

#### 5 Conclusion

In this survey, we provide a comprehensive overview of out-of-distribution generalization approaches from a causal perspective. In particular, depending on at which stage of training the machine learning framework the causal domain generalization component is applied, we classify them into three main categories, namely causal data augmentation methods (applied during the data pre-processing phase), invariant causal representation learning approaches (performed during the representation learning stage), and invariant causal mechanism learning algorithms (applied at the classifier level) and explain the state-of-the-art methods in each category. Depending on the approach taken, we further categorize the approaches in each category into sub-categories as shown in Figure 2. Moreover, we extend our categories to the textual and graph data and classify the approaches developed for those data types into our three categories. Comparing the comprehensive body of literature for the image data with recent works on textual and graph data, we observe many future research directions for these data types. Specifically, while most works on these data types belong to the causal data augmentation category, the causal representation learning and causal invariant mechanism learning-based approaches are greatly underexplored. We therefore suggest exploring these two directions for both of these data types. To evaluate the causal domain generalization approaches systematically, we provide a comprehensive list of commonly-used datasets and evaluation metrics to assess the performance of the proposed frameworks. These evaluation guidelines can be used by researchers and practitioners to appropriately evaluate the performance of their frameworks and compare the performance of existing methodologies.

#### References

- [1] Paras Sheth et al. Causal disentanglement with network information for debiased recommendations. *arXiv* preprint arXiv:2204.07221, 2022.
- [2] Jack Stilgoe. Machine learning, social learning and the governance of self-driving cars. *Social studies of science*, 2018.
- [3] Yann LeCun et al. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 1998.
- [4] Martin Arjovsky et al. Invariant risk minimization. arXiv preprint arXiv:1907.02893, 2019.
- [5] Rowland W Pettit, Robert Fullem, Chao Cheng, and Christopher I Amos. Artificial intelligence, machine learning, and deep learning for clinical outcome prediction. *Emerging topics in life sciences*, 2021.
- [6] Olga V Demler, Michael J Pencina, and Ralph B D'Agostino Sr. Impact of correlation on predictive ability of biomarkers. *Statistics in medicine*, 2013.
- [7] Kate Saenko et al. Adapting visual category models to new domains. In ECCV, 2010.
- [8] Zhihe Lu et al. Stochastic classifiers for unsupervised domain adaptation. In CVPR, 2020.
- [9] Ziwei Liu et al. Open compound domain adaptation. In CVPR, 2020.
- [10] Gilles Blanchard et al. Generalizing from several related classification tasks to a new unlabeled sample. *NeurIPS*, 2011.
- [11] Pierrick Bourrat. Measuring causal invariance formally. *Entropy*, 2021.
- [12] Peter Bühlmann. Invariance, causality and robustness. Statistical Science, 2020.
- [13] Divyat Mahajan et al. Domain generalization using causal matching. In ICML, 2021.
- [14] Kaiyang Zhou et al. Domain generalization: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- [15] Zheyan Shen et al. Towards out-of-distribution generalization: A survey. *arXiv preprint arXiv:2108.13624*, 2021.
- [16] Jindong Wang et al. Generalizing to unseen domains: A survey on domain generalization. IEEE TKDE, 2022.
- [17] Xiheng Zhang et al. Learning causal representation for training cross-domain pose estimator via generative interventions. In CVF, 2021.
- [18] Yuedong Chen et al. Towards unbiased visual emotion recognition via causal intervention. *arXiv preprint* arXiv:2107.12096, 2021.

- [19] Cheng Ouyang et al. Causality-inspired single-source domain generalization for medical image segmentation. *arXiv* preprint arXiv:2111.12525, 2021.
- [20] Jovana Mitrovic et al. Representation learning via invariant causal mechanisms. *arXiv preprint arXiv:2010.07922*, 2020.
- [21] Haoyue Bai et al. Decaug: Out-of-distribution generalization via decomposed feature representation and semantic augmentation. In *AAAI*, 2021.
- [22] Mathieu Chevalley, Charlotte Bunne, Andreas Krause, and Stefan Bauer. Invariant causal mechanisms through distribution matching. *arXiv* preprint arXiv:2206.11646, 2022.
- [23] Christina Heinze-Deml et al. Conditional variance penalties and domain shift robustness. *arXiv preprint arXiv:1710.11469*, 2017.
- [24] Chintan Trivedi et al. Contrastive learning of generalized game representations. In 2021 IEEE Conference on Games (CoG), 2021.
- [25] Chang Liu et al. Learning causal semantic representation for out-of-distribution prediction. NeurIPS, 2021.
- [26] Maggie Makar et al. Causally motivated shortcut removal using auxiliary labels. In AISTATS, 2022.
- [27] Jivat Neet Kaur et al. Modeling the data-generating process is necessary for out-of-distribution generalization. *arXiv e-prints*, 2022.
- [28] Judea Pearl. Causality. 2009.
- [29] Yang Chen et al. A style and semantic memory mechanism for domain generalization. In CVF, 2021.
- [30] Khurram Javed et al. Learning causal models online. arXiv preprint arXiv:2006.07461, 2020.
- [31] Mengyue Yang et al. Causalvae: Disentangled representation learning via neural structural causal models. In *CVF*, 2021.
- [32] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. arXiv preprint arXiv:1312.6114, 2013.
- [33] Yi-Fan Zhang et al. Learning domain invariant representations for generalizable person re-identification. *arXiv* preprint arXiv:2103.15890, 2021.
- [34] Xiang Deng and othersi. Comprehensive knowledge distillation with causal intervention. NeurIPS, 2021.
- [35] Tan Wang et al. Causal attention for unbiased visual recognition. In CVF, 2021.
- [36] Xin Li et al. Confounder identification-free causal visual feature learning. *arXiv preprint arXiv:2111.13420*, 2021
- [37] Kartik Ahuja et al. Invariance principle meets information bottleneck for out-of-distribution generalization. *NeurIPS*, 2021.
- [38] David Krueger et al. Out-of-distribution generalization via risk extrapolation (rex). In ICML, 2021.
- [39] Bo Li et al. Invariant information bottleneck for domain generalization. In AAAI, 2022.
- [40] Ruocheng Guo et al. Out-of-distribution prediction with invariant risk minimization: The limitation and an effective fix. *arXiv preprint arXiv:2101.07732*, 2021.
- [41] Dongsung Huh et al. The missing invariance principle found—the reciprocal twin of invariant risk minimization. *arXiv preprint arXiv:2205.14546*, 2022.
- [42] Alexis Bellot et al. Accounting for unobserved confounding in domain generalization. *arXiv preprint* arXiv:2007.10653, 2020.
- [43] Sreya Francis et al. Towards causal federated learning for enhanced robustness and privacy. *arXiv preprint arXiv:2104.06557*, 2021.
- [44] Jens Müller et al. Learning robust models using the principle of independent causal mechanisms. In *DAGM German Conference on Pattern Recognition*, 2021.
- [45] Arthur Gretton et al. Measuring statistical dependence with hilbert-schmidt norms. In *International conference on algorithmic learning theory*. Springer, 2005.
- [46] Xiangyu Zheng et al. Learning towards robustness in causally-invariant predictors. *arXiv preprint* arXiv:2107.01876, 2021.
- [47] Fangrui Lv et al. Causality inspired representation learning for domain generalization. In CVF, 2022.
- [48] Yunqi Wang et al. Contrastive ace: domain generalization through alignment of causal mechanisms. *arXiv* preprint arXiv:2106.00925, 2021.

- [49] Xinwei Sun, Botong Wu, Xiangyu Zheng, Chang Liu, Wei Chen, Tao Qin, and Tie-yan Liu. Latent causal invariant model. *arXiv preprint arXiv:2011.02203*, 2020.
- [50] Krikamol Muandet et al. Domain generalization via invariant feature representation. In ICML, 2013.
- [51] Dan Hendrycks et al. Pretrained transformers improve out-of-distribution robustness. *arXiv preprint* arXiv:2004.06100, 2020.
- [52] Zhao Wang et al. Identifying spurious correlations for robust text classification. *arXiv preprint arXiv:2010.02458*, 2020.
- [53] Ellery Wulczyn et al. Ex machina: Personal attacks seen at scale. In *Proceedings of the 26th international conference on world wide web*, 2017.
- [54] Seungtaek Choi et al. C21: Causally contrastive learning for robust text classification. 2022.
- [55] Zhao Wang et al. Robustness to spurious correlations in text classification via automatically generated counterfactuals. In AAAI, 2021.
- [56] Victor Veitch et al. Counterfactual invariance to spurious correlations: Why and how to pass stress tests. *arXiv* preprint arXiv:2106.00545, 2021.
- [57] Maxime Peyrard et al. Invariant language modeling. arXiv preprint arXiv:2110.08413, 2021.
- [58] Lu Cheng et al. Bias mitigation for toxicity detection via sequential decisions. In ACM SIGIR, 2022.
- [59] Amir Feder et al. Causal inference in natural language processing: Estimation, prediction, interpretation and beyond. *arXiv preprint arXiv:2109.00725*, 2021.
- [60] Megha Srivastava et al. Robustness to spurious correlations via human annotations. In ICML, 2020.
- [61] Divyansh Kaushik et al. Learning the difference that makes a difference with counterfactually-augmented data. *arXiv preprint arXiv:1909.12434*, 2019.
- [62] Fiona Anting Tan et al. Causal augmentation for causal sentence classification. In *Proceedings of the First Workshop on Causal Inference and NLP*, 2021.
- [63] Kartik Ahuja et al. Invariant risk minimization games. In ICML, 2020.
- [64] Gang Liu et al. Graph rationalization with environment-based augmentations. *arXiv preprint arXiv:2206.02886*, 2022.
- [65] Shaohua Fan et al. Generalizing graph neural networks on out-of-distribution graphs. *arXiv preprint arXiv:2111.10657*, 2021.
- [66] Yongduo Sui et al. Causal attention for interpretable and generalizable graph classification. *arXiv preprint* arXiv:2112.15089, 2022.
- [67] Beatrice Bevilacqua et al. Size-invariant graph representations for graph classification extrapolations. In *ICML*, 2021
- [68] Tom AB Snijders et al. Estimation and prediction for stochastic blockmodels for graphs with latent block structure. *Journal of classification*, 1997.
- [69] László Lovász et al. Limits of dense graph sequences. Journal of Combinatorial Theory, Series B, 2006.
- [70] Yongqiang Chen et al. Invariance principle meets out-of-distribution generalization on graphs. *arXiv* preprint *arXiv*:2202.05441, 2022.
- [71] Ying-Xin Wu et al. Discovering invariant rationales for graph neural networks. *arXiv preprint arXiv:2201.12872*, 2022.
- [72] Tero Karras et al. Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*, 2017.
- [73] Andreas Rossler et al. Faceforensics++: Learning to detect manipulated facial images. In *IEEE CVF*, 2019.
- [74] Liming Jiang et al. Deeperforensics-1.0: A large-scale dataset for real-world face forgery detection. In *IEEE CVF*, 2020.
- [75] Patrick Lucey et al. The extended cohn-kanade dataset (ck+): A complete dataset for action unit and emotion-specified expression. In *IEEE CVPR*, 2010.
- [76] Michel Valstar et al. Induced disgust, happiness and surprise: an addition to the mmi facial expression database. In *Proc. 3rd Intern. Workshop on EMOTION* (satellite of LREC): Corpora for Research on Emotion and Affect, 2010.

- [77] Guoying Zhao et al. Facial expression recognition from near-infrared videos. *Image and vision computing*, 2011.
- [78] Han Xiao et al. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv* preprint arXiv:1708.07747, 2017.
- [79] Yaroslav Ganin et al. Unsupervised domain adaptation by backpropagation. In ICML, 2015.
- [80] Yuval Netzer et al. Reading digits in natural images with unsupervised feature learning. 2011.
- [81] Muhammad Ghifary et al. Domain generalization for object recognition with multi-task autoencoders. In *IEEE CVF*, 2015.
- [82] Clifford R Jack Jr et al. The alzheimer's disease neuroimaging initiative (adni): Mri methods. *Journal of Magnetic Resonance Imaging: An Official Journal of the International Society for Magnetic Resonance in Medicine*, 2008.
- [83] Xiaosong Wang et al. Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. In CVPR, 2017.
- [84] Jeremy Irvin et al. Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison. In *AAAI*, 2019.
- [85] Tomás Franquet. Imaging of community-acquired pneumonia. Journal of thoracic imaging, 2018.
- [86] Catalin Ionescu et al. Human3. 6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE transactions on pattern analysis and machine intelligence*, 2013.
- [87] Timo Von Marcard et al. Recovering accurate 3d human pose in the wild using imus and a moving camera. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.
- [88] Dushyant Mehta et al. Monocular 3d human pose estimation in the wild using improved cnn supervision. In 2017 international conference on 3D vision (3DV), 2017.
- [89] Gul Varol et al. Learning from synthetic humans. In CVPR, 2017.
- [90] Leonid Sigal et al. Humaneva: Synchronized video and motion capture dataset and baseline algorithm for evaluation of articulated human motion. *International journal of computer vision*, 2010.
- [91] Wei Li et al. Locally aligned feature transforms across views. In CVPR, 2013.
- [92] Wei Li et al. Deepreid: Deep filter pairing neural network for person re-identification. In CVPR, 2014.
- [93] Liang Zheng et al. Scalable person re-identification: A benchmark. In CVPR, 2015.
- [94] Zhedong Zheng et al. Unlabeled samples generated by gan improve the person re-identification baseline in vitro. In CVPR, 2017.
- [95] Tong Xiao et al. End-to-end deep learning for person search. arXiv preprint arXiv:1604.01850, 2016.
- [96] Catherine Wah et al. The caltech-ucsd birds-200-2011 dataset. 2011.
- [97] Bolei Zhou et al. Places: A 10 million image database for scene recognition. *IEEE transactions on pattern analysis and machine intelligence*, 2017.
- [98] Sara Beery et al. The iwildcam 2018 challenge dataset. arXiv preprint arXiv:1904.05986, 2019.
- [99] Sara Beery et al. Recognition in terra incognita. In ECCV, 2018.
- [100] Da Li et al. Deeper, broader and artier domain generalization. In *ICCV*, 2017.
- [101] Hemanth Venkateswara et al. Deep hashing network for unsupervised domain adaptation. In CVPR, 2017.
- [102] Xingchao Peng et al. Moment matching for multi-source domain adaptation. In IEEE CVF, 2019.
- [103] Dan Hendrycks et al. The many faces of robustness: A critical analysis of out-of-distribution generalization. In *IEEE CVF*, 2021.
- [104] Chen Fang et al. Unbiased metric learning: On the utilization of multiple datasets and web images for softening bias. In *ICCV*, 2013.
- [105] Dan Hendrycks et al. Benchmarking neural network robustness to common corruptions and perturbations. *arXiv* preprint arXiv:1903.12261, 2019.
- [106] Yue and others He. Towards non-iid image classification: A dataset and baselines. Pattern Recognition, 2021.
- [107] Xingxuan Zhang et al. Nico++: Towards better benchmarking for domain generalization. *arXiv preprint* arXiv:2204.08040, 2022.
- [108] Zhenqin Wu et al. Moleculenet: a benchmark for molecular machine learning. Chemical science, 2018.

- [109] Rafael Gómez-Bombarelli et al. Automatic chemical design using a data-driven continuous representation of molecules. *ACS central science*, 2018.
- [110] Aleksandar Bojchevski et al. Deep gaussian embedding of graphs: Unsupervised inductive learning via ranking. *arXiv preprint arXiv:1707.03815*, 2017.
- [111] Weihua Hu et al. Open graph benchmark: Datasets for machine learning on graphs. NeurIPS, 2020.
- [112] Zhitao Ying et al. Gnnexplainer: Generating explanations for graph neural networks. NeurIPS, 2019.
- [113] Fang Fang et al. Domain adaptation for sentiment classification in light of multiple sources. *INFORMS Journal on Computing*, 2014.
- [114] Andrew Maas et al. Learning word vectors for sentiment analysis. In ACL: Human language technologies, 2011.
- [115] Julian John McAuley et al. From amateurs to connoisseurs: modeling the evolution of user expertise through online reviews. In WWW, 2013.
- [116] Daniel Borkan et al. Nuanced metrics for measuring unintended bias with real data for text classification. In *WWW*, 2019.
- [117] Adina Williams et al. A broad-coverage challenge corpus for sentence understanding through inference. *arXiv* preprint arXiv:1704.05426, 2017.
- [118] Chaochao Lu et al. Invariant causal representation learning for out-of-distribution generalization. In ICLR, 2021.
- [119] Ruoyu Wang et al. Improving ood generalization with causal invariant transformations. 2021.
- [120] Justin B Kinney et al. Equitability, mutual information, and the maximal information coefficient. *Proceedings of the National Academy of Sciences*, 2014.