

# Understanding Hessian Alignment for Domain Generalization

Sobhan Hemati\*      Guojun Zhang\*      Amir Estiri      Xi Chen  
Huawei Noah’s Ark Lab

{sobhan.hemati, guojun.zhang, amir.hossein.estiri1, xi.chen4}@huawei.com

## Abstract

*Out-of-distribution (OOD) generalization is a critical ability for deep learning models in many real-world scenarios including healthcare and autonomous vehicles. Recently, different techniques have been proposed to improve OOD generalization. Among these methods, gradient-based regularizers have shown promising performance compared with other competitors. Despite this success, our understanding of the role of Hessian and gradient alignment in domain generalization is still limited. To address this shortcoming, we analyze the role of the classifier’s head Hessian matrix and gradient in domain generalization using recent OOD theory of transferability. Theoretically, we show that spectral norm between the classifier’s head Hessian matrices across domains is an upper bound of the transfer measure, a notion of distance between target and source domains. Furthermore, we analyze all the attributes that get aligned when we encourage similarity between Hessians and gradients. Our analysis explains the success of many regularizers like CORAL, IRM, V-REx, Fish, IGA, and Fishr as they regularize part of the classifier’s head Hessian and/or gradient. Finally, we propose two simple yet effective methods to match the classifier’s head Hessians and gradients in an efficient way, based on the Hessian Gradient Product (HGP) and Hutchinson’s method (Hutchinson), and without directly calculating Hessians. We validate the OOD generalization ability of proposed methods in different scenarios, including transferability, severe correlation shift, label shift and diversity shift. Our results show that Hessian alignment methods achieve promising performance on various OOD benchmarks. The code is available [here](#).*

## 1. Introduction

A typical assumption in the design of current supervised deep learning algorithms is identical distributions of test and training data. Unfortunately, in real-world problems, this assumption is not always true (Koyama and Yamaguchi,

2021) and the distribution gap between test and training data degrades the performance of deep learning models (Wang et al., 2022). In practice, not only is there a distribution shift between test and training data, but also the training data does not follow the i.i.d. assumption but contains data from multiple domains. Histopathology medical datasets are a concrete example where each hospital is using different types of staining procedures and/or scanners (Chang et al., 2021). In such a scenario, a deep learning model often fails to generalize its knowledge to unseen domains. Arjovsky et al. (2019) argue that the main reason behind poor OOD generalization is the tendency of deep learning models to capture simple spurious correlations instead of real causal information. In order to improve OOD generalization, the learning algorithm has to learn from invariant mechanisms.

The problem setup in domain generalization (DG) assumes the training set contains data from multiple sources (domains) where the task remains the same across different data sources (Muandet et al., 2013). Furthermore, we assume there is a causal mechanism that remains invariant across different domains and hopefully for OOD data. Although many learning algorithms have been proposed to capture invariance, a recent work by Gulrajani and Lopez-Paz (2020) shows that given a standard experimental setting, the classic Empirical Risk Minimization (ERM) algorithm (Vapnik, 1999), which minimizes the average of losses across different training domains, outperforms many recently proposed domain generalization methods. This suggests that many current algorithms may not be successful in capturing the invariance in data.

Out of the effective DG algorithms, one recent and promising line of research is gradient-based methods which try to capture invariance in gradient space (Parascandolo et al., 2020; Koyama and Yamaguchi, 2021; Shi et al., 2021; Rame et al., 2022). One fast and simple explanation for this success can be derived from the seminal work by Jaakkola and Haussler (1998). According to this work, given a generative model, the gradient with respect to each parameter quantifies the role of that parameter in the generation of a data point. With this intuition, encouraging the similarity between gradients across environments can be translated as

\*Equal contribution. ICCV 2023 camera ready.

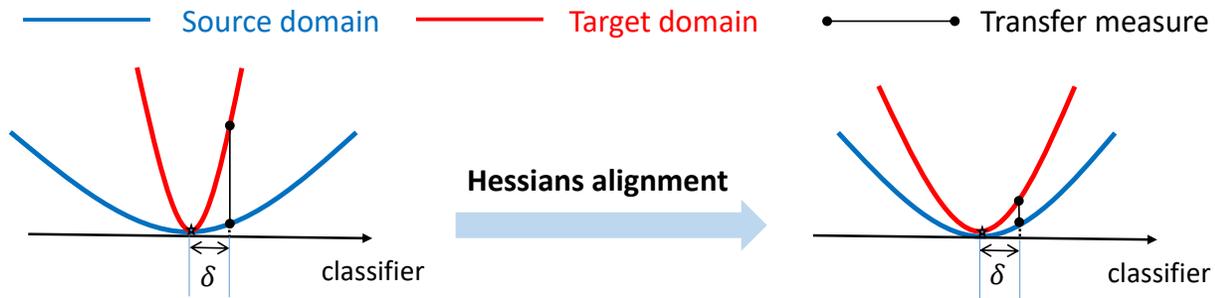


Figure 1: Visualization of Hessian alignment for domain generalization.  $\delta$  denotes the small deviation from an optimum. Hessian alignment matches the curvature and thus improves the transfer measure between source and target domains.

enforcing the network parameters to contribute the same for different environments which captures a notion of invariance. Beyond gradient matching, the Invariant Learning Consistency (ILC) measure (Parascandolo et al., 2020) motivates matching Hessians across different environments after convergence, given that we found the optimal solution for all environments and that the Hessian matrix is diagonal. Although the efforts in Parascandolo et al. (2020); Koyama and Yamaguchi (2021); Rame et al. (2022) improve our understanding of the role of Hessians and gradients to some extent, the ILC measure is built based on heuristics and is only valid under restricted assumptions like diagonal Hessian matrices. Moreover, it is not exactly clear what attributes are getting aligned when we match gradients or even Hessians. Finally, the proposed methods to align gradients or Hessians (e.g., ANDmask, Parascandolo et al. 2020) seem to be suboptimal as discussed by follow-up works (Shahtalebi et al., 2021; Rame et al., 2022).

To address these limitations, in this paper, we study the role of Hessians and gradients in domain generalization. Since Hessians for the whole neural networks are hard to compute, we focus on Hessians of the classifier head, which contains more information than one usually expects. We summarize our contributions as follows:

- To justify Hessian matching, we utilize a recent concept from statistical learning theory, called transfer measure (Zhang et al., 2021), which describes transferability between target and source domains. Unlike ILC, transfer measure avoids the restrictive assumptions of diagonal Hessians. We theoretically show that the distance between the classifier’s head Hessians is an upper bound of the transfer measure.
- Furthermore, we show that Hessians and gradient alignment can be treated as feature matching. Our analysis of feature matching compares other DG algorithms like CORAL (Sun and Saenko, 2016) and V-REx (Krueger

et al., 2021) in a unified framework. Especially, the success of CORAL can be attributed to approximate Hessian alignment using our analysis.

- Last but not least, to match Hessians efficiently, we propose two simple yet effective methods, based on different estimation of the Hessian. To our knowledge, these methods are the first DG algorithms based on Hessian estimation that align Hessians and gradients simultaneously.

Figure 1 provides an intuitive explanation why Hessian alignment can reduce domain shift through minimizing transferability. With Hessian matching, the transfer measure, or the excess risk gap between source and target domains, is minimized. This makes any near-optimal source classifier to be also near-optimal on the target, and thus improves the transferability and OOD generalization.

## 2. Problem Setting and Related Work

Consider a deep neural network that consists of a feature extractor  $g$  and a classifier head  $h_\theta$  with parameters  $\theta$ , and suppose the training data contains  $n$  different source domains  $\mathcal{E}^{source} = \{\mathcal{S}_1, \dots, \mathcal{S}_n\}$ . The main goal in DG is to design a learning algorithm where the trained model can generalize to an unseen target domain  $\mathcal{T}$ . Denoting  $\mathcal{H}$  as the hypothesis class, the classification loss of a classifier  $h_\theta \in \mathcal{H}$  on a domain  $\mathcal{D}$  is:

$$\mathcal{L}_{\mathcal{D}}(\theta) = \mathbb{E}_{(x,y) \sim \mathcal{D}}[\ell(\hat{y}, y; \theta)], \quad (1)$$

where  $x$  and  $y$  are the input and the associated one-hot label,  $\hat{y}$  is the logit (classifier output) and  $\hat{y} = h_\theta(g(x))$  where  $g(x)$  is the output of a feature extractor,  $h_\theta$  is the classifier and  $\ell(\hat{y}, y; \theta)$  is the cross entropy between  $\hat{y}$  and  $y$ . The classic baseline for OOD generalization is Empirical Risk Minimization (ERM, Vapnik, 1999) which minimizes the average of losses across different training environments

namely  $\frac{1}{n} \sum_e \mathcal{L}_{S_e}$ . However, it has been shown that this approach might fail to generalize to out-of-distribution data when there is an easy-to-learn spurious correlation in data (Arjovsky et al., 2019). To tackle this, many algorithms are proposed to avoid domain-specific representations and capture invariance of data from different perspectives. We categorize the most related into four main classes:

**Data augmentation.** Recently many papers have shown that data augmentation can improve OOD generalization. Gulrajani and Lopez-Paz (2020) showed that the classic ERM with data augmentation can outperform many DG algorithms. Zhang et al. (2018) proposed *mixup*, a data augmentation strategy designed for domain generalization where samples across multiple source domains get mixed to generate new data for training. Ilse et al. (2021) derived an algorithm that select proper data augmentation.

**Robust optimization** techniques have been used to achieve better OOD generalization (Sagawa et al., 2020; Hu et al., 2018). These techniques minimize the worst-case loss over a set of source domains, and can be helpful when we know the distance between training and test environments.

**Learning domain invariant features.** Invariant representation learning was initialized by seminal work of Ben-David et al. (2010). A wide range of notions of invariance have been proposed since then. Tzeng et al. (2014) proposed feature matching across domains while Sun and Saenko (2016) considered encouraging the similarity between feature covariance matrices. DANN (Ganin et al., 2016) proposed matching feature distributions with an adversarial network, which is followed by MMD (Li et al., 2018b) and MDAN (Zhao et al., 2018). Zhao et al. (2019) discussed the fundamental tradeoff between feature matching and optimal joint risk. Another line work started from Invariant Risk Minimization (IRM) (Arjovsky et al., 2019), which proposed to encourage the classifier head to be optimal for different domains simultaneously. This idea led to a new notion of invariance as followed by e.g., Risk Extrapolation (V-REx) (Krueger et al., 2021) where authors proposed similarity between losses across training domains.

**Gradient matching based algorithms.** Recently, gradient alignment has been used for domain generalization (Parascandolo et al., 2020; Shahtalebi et al., 2021; Koyama and Yamaguchi, 2021; Mansilla et al., 2021; Shi et al., 2021; Rame et al., 2022) with good performance. In this regard, Shi et al. (2021) proposed Fish to maximize the inner product of average gradients for different domains. In a similar approach, Koyama and Yamaguchi (2021) proposed the Inter Gradient Alignment (IGA) algorithm that minimizes the variance of average gradients over environments. Parascandolo et al. (2020) developed ANDmask algorithm where the authors proposed ILC, which measures the consistency of local minima for different environments. Parascandolo et al. (2020) proposed ANDmask to reduce the speed of

convergence in directions where landscapes have different curvatures. Inspired by ILC measure (Parascandolo et al., 2020), Rame et al. (2022) proposed Fishr as a better alternative compared with ANDmask to reduce inconsistency. In Fishr, the variance of domain level gradients is minimized.

### 3. Theory of Hessian Matching

In this section, we study the role of Hessians and gradients in domain generalization from two views: transferability and invariant representation learning.

#### 3.1. Hessian alignment minimizes transfer measure

We motivate our work through the lens of transferability and transfer measure introduced in Zhang et al. (2021). By their definition, domain  $\mathcal{S}$  is transferable to domain  $\mathcal{T}$  on a hypothesis class  $\mathcal{H}$ , if given any near-optimal classifier  $h_\theta$  on  $\mathcal{S}$ , it also achieves near-optimal performance on  $\mathcal{T}$ . More precisely, we have:

**Definition 1 (transferability, Zhang et al. 2021).**  $\mathcal{S}$  is  $(\delta_S, \delta_T)_{\mathcal{H}}$ -transferable to  $\mathcal{T}$  if for  $\delta_S > 0$ , there exists  $\delta_T > 0$  such that  $\text{argmin}(\mathcal{L}_{\mathcal{S}}, \delta_S)_{\mathcal{H}} \subseteq \text{argmin}(\mathcal{L}_{\mathcal{T}}, \delta_T)_{\mathcal{H}}$ , where:

$$\begin{aligned} \text{argmin}(\mathcal{L}_{\mathcal{D}}, \delta_{\mathcal{D}})_{\mathcal{H}} &:= \{h_\theta \in \mathcal{H} : \mathcal{L}_{\mathcal{D}}(\theta)\} \\ &\leq \inf_{h_\theta \in \mathcal{H}} \mathcal{L}_{\mathcal{D}}(\theta) + \delta_{\mathcal{D}}, \mathcal{D} \in \{\mathcal{S}, \mathcal{T}\}. \end{aligned}$$

The set  $\text{argmin}(\mathcal{L}_{\mathcal{D}}, \delta_{\mathcal{D}})_{\mathcal{H}}$  is called a  $\delta$ -minimal set (Koltchinskii, 2010) of  $\mathcal{L}_{\mathcal{D}}$ , and represents the near-optimal set of classifiers. With the  $\delta$ -minimal set, we can define the transfer measure:

**Definition 2 (transfer measure, Zhang et al. 2021).** Given  $\Gamma = \text{argmin}(\mathcal{L}_{\mathcal{S}}, \delta_S)_{\mathcal{H}}$ ,  $\mathcal{L}_{\mathcal{S}}^* := \inf_{h_\theta \in \Gamma} \mathcal{L}_{\mathcal{S}}(\theta)$  and  $\mathcal{L}_{\mathcal{T}}^* := \inf_{h_\theta \in \Gamma} \mathcal{L}_{\mathcal{T}}(\theta)$  we define the transfer measure  $T_\Gamma(\mathcal{S} \parallel \mathcal{T})$  as:

$$T_\Gamma(\mathcal{S} \parallel \mathcal{T}) := \sup_{h_\theta \in \Gamma} \mathcal{L}_{\mathcal{T}}(\theta) - \mathcal{L}_{\mathcal{T}}^* - (\mathcal{L}_{\mathcal{S}}(\theta) - \mathcal{L}_{\mathcal{S}}^*). \quad (2)$$

The transfer measure in Eq. 2 measures the upper bound of the difference between source domain excess risk ( $\mathcal{L}_{\mathcal{S}}(\theta) - \mathcal{L}_{\mathcal{S}}^*$ ) and the target domain excess risk ( $\mathcal{L}_{\mathcal{T}}(\theta) - \mathcal{L}_{\mathcal{T}}^*$ ), on a neighborhood  $\Gamma$ . Usually, we choose  $\Gamma$  to be a neighborhood of our learned near-optimal source classifier.

Now, we state our main theorem which, under mild assumptions, implies that Hessian alignment effectively reduces the transfer measure and improving transferability.

**Theorem 3 (Hessian alignment).** Suppose both the source and target domain losses  $\mathcal{L}_{\mathcal{S}}$  and  $\mathcal{L}_{\mathcal{T}}$  are  $\mu$ -strongly convex with respect to the classifier head and twice differentiable with the same minimizer, i.e.,  $\text{argmin} \mathcal{L}_{\mathcal{S}} = \text{argmin} \mathcal{L}_{\mathcal{T}}$ . Then with  $\delta = 2\delta_S/\mu$  and  $\Gamma = \text{argmin}(\mathcal{L}_{\mathcal{S}}, \delta_S)_{\mathcal{H}}$  we have:

$$T_\Gamma(\mathcal{S} \parallel \mathcal{T}) \leq \frac{1}{2} \delta^2 \|\mathbf{H}_{\mathcal{T}} - \mathbf{H}_{\mathcal{S}}\|_2 + o(\delta^2). \quad (3)$$

*Proof.* Let  $\theta^*$  represent the optimal classifier for both source and target domains. Given that  $\mathcal{L}_S(\theta)$  is strongly convex, we can write  $\inf_{h_{\theta} \in \mathcal{H}} \mathcal{L}_S(\theta) = \mathcal{L}_S(\theta^*)$  and subsequently  $\Gamma$  can be written as

$$\operatorname{argmin}(\mathcal{L}_S, \delta_S)_{\mathcal{H}} := \{h_{\theta} \in \mathcal{H} : \mathcal{L}_S(\theta) - \mathcal{L}_S(\theta^*) \leq \delta_S\}.$$

Now, we define set  $\mathcal{F}_1$  as

$$\mathcal{F}_1 = \{\theta : \mathcal{L}_S(\theta) - \mathcal{L}_S(\theta^*) \leq \delta_S\}. \quad (4)$$

On the other hand, from the  $\mu$ -strong convexity of  $\mathcal{L}_S(\theta)$ , we can write

$$\mathcal{L}_S(\theta) \geq \mathcal{L}_S(\theta^*) + \nabla_{\theta} \mathcal{L}_S^{\top}(\theta^*)(\theta - \theta^*) + \frac{\mu}{2} \|\theta - \theta^*\|_2^2. \quad (5)$$

Given the optimality of  $\theta^*$  (and thus  $\nabla_{\theta} \mathcal{L}_S(\theta^*) = \mathbf{0}$ ), from Eqs. 4 and 5 we obtain that for any  $\theta \in \mathcal{F}_1$ ,

$$\frac{\mu}{2} \|\theta - \theta^*\|_2^2 \leq \mathcal{L}_S(\theta) - \mathcal{L}_S(\theta^*) \leq \delta_S. \quad (6)$$

If we define set  $\mathcal{F}_2$  as

$$\mathcal{F}_2 = \{\theta : \frac{\mu}{2} \|\theta - \theta^*\|_2^2 \leq \delta_S\}. \quad (7)$$

Then (6) implies  $\mathcal{F}_1 \subset \mathcal{F}_2$  and as a result, an upper bound of the transfer measure  $T_{\Gamma}(S||\mathcal{T})$  can be written as

$$\sup_{\|\theta - \theta^*\|_2 \leq \delta} \mathcal{L}_{\mathcal{T}}(\theta) - \mathcal{L}_{\mathcal{T}}(\theta^*) - (\mathcal{L}_S(\theta) - \mathcal{L}_S(\theta^*)), \quad (8)$$

with  $\delta = 2\delta_S/\mu$ . Now, we may write the second-order Taylor expansion  $\mathcal{L}_S$  around  $\theta^*$  as:

$$\begin{aligned} \mathcal{L}_S(\theta) &= \mathcal{L}_S(\theta^*) + (\theta - \theta^*)^{\top} \nabla_{\theta} \mathcal{L}_S(\theta^*) + \\ &+ \frac{1}{2} (\theta - \theta^*)^{\top} \mathbf{H}_S (\theta - \theta^*) + o(\|\theta - \theta^*\|^2), \end{aligned} \quad (9)$$

where given  $\nabla_{\theta} \mathcal{L}_S(\theta^*) = \mathbf{0}$ , the source domain excess risk is

$$\mathcal{L}_S(\theta) - \mathcal{L}_S(\theta^*) = \frac{1}{2} \Delta \theta^{\top} \mathbf{H}_S \Delta \theta + o(\|\theta - \theta^*\|^2), \quad (10)$$

with  $\Delta \theta = \theta - \theta^*$ . Following a similar path for the target domain and considering that  $\theta^*$  is the optimal solution for both source and target domains, the transfer measure  $T_{\Gamma}(S||\mathcal{T})$  is upper bounded by

$$\sup_{\|\theta - \theta^*\|_2 \leq \delta} \frac{1}{2} (\theta - \theta^*)^{\top} (\mathbf{H}_{\mathcal{T}} - \mathbf{H}_S) (\theta - \theta^*) + o(\delta^2), \quad (11)$$

which, from the definition of spectral norm, results in

$$T_{\Gamma}(S||\mathcal{T}) \leq \frac{1}{2} \delta^2 \|\mathbf{H}_{\mathcal{T}} - \mathbf{H}_S\|_2 + o(\delta^2), \quad (12)$$

and the proof is complete.  $\square$

Note that convexity and differentiability are easily satisfied if we use the standard linear classifier head and cross entropy loss with softmax. Since in practice we do not have access to the target domain, we minimize the distance between available source domains. For simplicity, we replace the spectral norm with the Frobenius norm (note that Frobenius norm is an upper bound, i.e.,  $\|\cdot\|_2 \leq \|\cdot\|_F$ ). As opposed to pairwise regularizers which grow quadratically with the number of environments, we minimize the variance of Hessian distances across domains, namely

$$\frac{1}{n} \sum_{e=1}^n \|\mathbf{H}_{S_e} - \overline{\mathbf{H}}_S\|_2^2 \text{ where } \overline{\mathbf{H}}_S = \frac{1}{n} \sum_{e=1}^n \mathbf{H}_{S_e}$$

and  $\mathbf{H}_{S_e}$  is the Hessian matrix for the  $e$ -th source domain. In this way, computation is linear with respect to the number of environments.

### 3.2. Hessian and gradient alignment is feature matching

In §3.1 we assumed that the minimizers of source and target are the same. In order to align them, gradient alignment is often necessary. To get more insight into the attributes that get matched while the Hessians and gradients get aligned, in Proposition 4, we study the Hessian and gradient structures with respect to the classifier head parameters.

**Proposition 4 (Alignment attributes in Hessian and gradient).** *Let  $\hat{y}_p$  and  $y_p$  be the network prediction and true target with the  $p$ -th class,  $z_i$  be the  $i$ -th feature before the classifier. Denote the classifier's parameters as  $w_{k,i}$ , the element in row  $k$  and column  $i$  of the classifier weight matrix, and  $b_k$  as the bias term for the  $k$ -th neuron. Matching the gradients and Hessians w.r.t. the classifier head across domains aligns the following:*

$$\frac{\partial \ell}{\partial b_p} = (\hat{y}_p - y_p), \text{ (Error)} \quad (13)$$

$$\frac{\partial \mathcal{L}}{\partial w_{p,q}} = (\hat{y}_p - y_p) z_q, \text{ (Error} \times \text{Feature)} \quad (14)$$

$$\frac{\partial^2 \ell}{\partial b_u \partial b_v} = \hat{y}_u (\delta_{u,v} - \hat{y}_v), \text{ (Logit)} \quad (15)$$

$$\frac{\partial^2 \ell}{\partial w_{p,q} \partial b_u} = z_q \hat{y}_p (\delta_{p,u} - \hat{y}_u), \text{ (Logit} \times \text{Feature)} \quad (16)$$

$$\frac{\partial^2 \ell}{\partial w_{p,q} \partial w_{u,v}} = \hat{y}_p z_q z_v (\delta_{p,u} - \hat{y}_u), \text{ (Logit} \times \text{Covariance)} \quad (17)$$

where  $\delta_{p,u}$  is the Kronecker delta which is 1 if  $p = u$  and 0 otherwise and  $u, v, p, q, k, i$  are dummy indices.

Here the error means the difference between the true one-hot label and our prediction. Eqs. 13–17 show that

Table 1: The alignment attributes for different Domain Generalization algorithms where Loss is  $(y_p - \hat{y}_p)^2$ , Feature  $z_q$ , Covariance  $z_q z_v$ , Error  $y_p - \hat{y}_p$ , Error  $\times$  Feature  $(\hat{y}_p - y_p)z_q$ , Logit  $\hat{y}_u(\delta_{u,v} - \hat{y}_v)$ , Logit  $\times$  Feature  $z_q \hat{y}_p(\delta_{p,u} - \hat{y}_u)$ , and Logit  $\times$  Covariance is  $\hat{y}_p z_q z_v(\delta_{p,u} - \hat{y}_u)$ . Gradient alignment method matches Error and Error  $\times$  Feature while matching Hessians would align Logit, Logit  $\times$  Feature, Logit  $\times$  Covariance. <sup>1</sup>To be precise, the Loss that get aligned in V-Rex is negative log likelihood instead of squared error. <sup>2</sup>To be precise, Fishr aligns Loss  $\times$  Feature.

Alignment attribute	Loss	Feature	Covariance	Error	Error $\times$ Feature	Logit	Logit $\times$ Feature	Logit $\times$ Covariance
V-Rex	✓ <sup>1</sup>	✗	✗	✗	✗	✗	✗	✗
CORAL	✗	✓	✓	✗	✗	✗	✗	✗
IGA	✗	✗	✗	✓	✓	✗	✗	✗
Fish	✗	✗	✗	✓	✓	✗	✗	✗
Fishr	✓	✗	✗	✗	✓ <sup>2</sup>	✗	✗	✗
Hessian Alignment	✗	✗	✗	✓	✓	✓	✓	✓

matching Hessians and gradients with respect to classifier parameters will align the errors, features weighted by errors, logits, features weighted by logits, and covariance weighted by logits across different domains simultaneously. Although gradient alignment can be helpful before convergence, when the training is close to convergence, the gradients go to zero and aligning them is not helpful anymore. On the other hand, the attributes extracted from the Hessian matrix can remain non-zero both before and after convergence and aligning them can boost OOD performance.

Proposition 4 shows that matching gradients and Hessians can be seen as a generalization of other works like V-Rex (Krueger et al., 2021), CORAL (Sun and Saenko, 2016), Fish (Shi et al., 2021), IGA (Koyama and Yamaguchi, 2021), and Fishr (Rame et al., 2022), as these algorithms only partially realize the matching of loss/error, feature or covariance. We summarize alignments in different methods in Table 1 and reveal their connection to Hessian and gradient alignment. One interesting observation is that similar to Hessian alignment, CORAL also matches the feature and its covariance, which opens the venue to understand the success of CORAL with Hessian matching. In the supplementary, we also present similar results for regression tasks.

## 4. Efficient Hessian matching

So far, we have shown that aligning Hessians across domains can reduce the transfer measure (or in other words, increase transferability) and match the representations at different levels i.e., logits, features, errors and covariances. However, given the computational complexity of calculating Hessian matrices, it is quite challenging to directly minimize  $\|\mathbf{H}_{S_1} - \mathbf{H}_{S_2}\|_F$ . To align Hessian matrices efficiently, we propose two approaches, aligning Hessian-gradient products and matching Hessian diagonals using Hutchinson’s trace estimator (Bekas et al., 2007).

### 4.1. Hessian Gradient Product (HGP)

For the Hessian-gradient product matching we propose the following loss:

$$\mathcal{L}_{\text{HGP}} = \frac{1}{n} \sum_{e=1}^n \mathcal{L}_{S_e} + \alpha \|\mathbf{H}_{S_e} \nabla_{\theta} \mathcal{L}_{S_e} - \overline{\mathbf{H}_S \nabla_{\theta} \mathcal{L}_S}\|_2^2 + \beta \|\nabla_{\theta} \mathcal{L}_{S_e} - \overline{\nabla_{\theta} \mathcal{L}_S}\|_2^2, \quad (18)$$

where  $\alpha$  and  $\beta$  are regularization parameters,

$$\overline{\nabla_{\theta} \mathcal{L}_S} = \frac{1}{n} \sum_{e=1}^n \nabla_{\theta} \mathcal{L}_{S_e}, \quad \text{and} \quad (19)$$

$$\overline{\mathbf{H}_S \nabla_{\theta} \mathcal{L}_S} = \frac{1}{n} \sum_{e=1}^n \mathbf{H}_{S_e} \nabla_{\theta} \mathcal{L}_{S_e}. \quad (20)$$

Given that  $\nabla_{\theta} \mathcal{L}_{S_e}$  and  $\overline{\nabla_{\theta} \mathcal{L}_S}$  are close enough, minimizing eq. 18 reduces the distance between  $\mathbf{H}_{S_e}$  and  $\overline{\mathbf{H}_S}$ . We can efficiently compute the HGP term without directly calculating Hessians, using the following expression

$$\mathbf{H}_{S_e} \nabla_{\theta} \mathcal{L}_{S_e} = \|\nabla_{\theta} \mathcal{L}_{S_e}\| \cdot \nabla_{\theta} \|\nabla_{\theta} \mathcal{L}_{S_e}\|, \quad (21)$$

which follows from  $\nabla_{\mathbf{x}} \|\mathbf{x}\|_2 = \mathbf{x} / \|\mathbf{x}\|_2$  with  $\mathbf{x} \in \mathbb{R}^d$  and chain rule. For domain  $e$ , to calculate the exact value of the Hessian-gradient product, we only need two rounds of back-propagation: one for  $\nabla_{\theta} \mathcal{L}_{S_e}$  and another for  $\nabla_{\theta} \|\nabla_{\theta} \mathcal{L}_{S_e}\|$ .

### 4.2. Hutchinson’s method

In order to estimate Hessians more accurately, we propose another approach to estimate the Hessian diagonal using Hutchinson’s method (Bekas et al., 2007):

$$\mathbf{D}_{S_e} = \text{diag}(\mathbf{H}_{S_e}) = \mathbb{E}_{\mathbf{r}}[\mathbf{r} \odot (\mathbf{H}_{S_e} \cdot \mathbf{r})], \quad (22)$$

where  $\mathbb{E}_{\mathbf{r}}$  is expectation of random variable  $\mathbf{r}$  sampled from Rademacher distribution and  $\odot$  is Hadamard product. In

practice, the expectation is replaced with the mean over a set of samples. Note that  $\mathbf{H}_{S_e} \cdot \mathbf{r}$  can be calculated efficiently by calculating  $\nabla_{\theta} (\nabla_{\theta} \mathcal{L}_{S_e}^{\top} \mathbf{r})$ . Given that we have estimated the Hessian diagonal, the  $\mathcal{L}_{\text{Hutchinson}}$  loss is

$$\mathcal{L}_{\text{Hutchinson}} = \frac{1}{n} \sum_{e=1}^n \mathcal{L}_{S_e} + \alpha \|\mathbf{D}_{S_e} - \overline{\mathbf{D}_S}\|_2^2 + \beta \|\nabla_{\theta} \mathcal{L}_{S_e} - \overline{\nabla_{\theta} \mathcal{L}_S}\|_2^2. \quad (23)$$

Compared to HGP, Hutchinson’s method is more expensive to compute, as it requires an estimation with sampling. In practice, we use 100 samples for each domain. The trade-off is that Hutchinson’s method gives us a more accurate matching algorithm for Hessian, since after training, the Hessian matrix is often close to diagonal (Skorski, 2021), whereas the gradient (and thus the Hessian-gradient product) becomes nearly zero.

## 5. Experiments

To comprehensively validate the effectiveness of HGP and Hutchinson, we evaluate the proposed algorithms from multiple views, in terms of transferability (Zhang et al., 2021), OOD generalization on datasets with correlation shift (Colored MNIST), correlation shift with label shift (imbalanced Colored MNIST), and diversity shift (Ye et al., 2021), including PACS (Li et al., 2017) and OfficeHome (Venkateswara et al., 2017) datasets from DomainBed.

### 5.1. Transferability

To show the role of Hessian discrepancy in improving transferability (i.e., Theorem 3), we use the algorithm presented in Zhang et al. (2021), which computes the worst-case gap  $\sup_{\|\theta - \theta^*\| \leq \delta} \mathcal{L}_{S_i}(\theta) - \mathcal{L}_{S_j}(\theta)$  among all pairs of  $(i, j)$ , to evaluate transferability among multiple domains. First, we find the domains  $i$  and  $j$  such that  $\mathcal{L}_{S_i}$  and  $\mathcal{L}_{S_j}$  are maximized and minimized respectively and consider the  $\mathcal{L}_{S_i} - \mathcal{L}_{S_j}$  as the gap. Then, we run an ascent optimizer on classifier  $\theta$  to maximize gap and project  $\theta$  on to Euclidean ball  $\|\theta - \theta^*\| \leq \delta$ . We repeat this procedure for multiple rounds and report the test target accuracy for the worst-case gap. Figure 2 shows that OOD accuracies of both HGP and Hutchinson are robust against the attack of an ascent optimizer, which implies that minimizing  $\|\mathbf{H}_{S_1} - \mathbf{H}_{S_2}\|_F$  improves transferability.

### 5.2. OOD generalization under correlation shift using Colored MNIST

Our Colored MNIST (CMNIST) setup is from Arjovsky et al. (2019), where the dataset contains two training and one test environments and the objective is a binary classification task in which class zero and one contain digits less than 5 and greater and equal 5 respectively. Further, the labels have

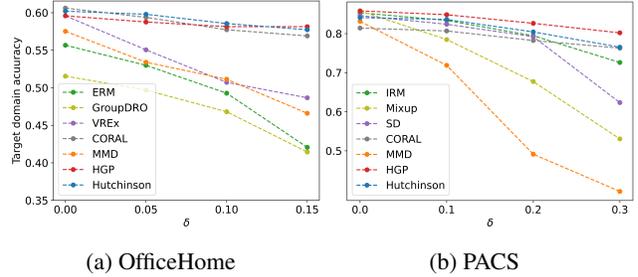


Figure 2: OOD accuracy of HGP and Hutchinson compared with multiple baselines on OfficeHome and PACS datasets with different values of  $\delta$  for the ascent optimizer.

been flipped with a probability 0.25 and each image has been colored either red or yellow such that the colors are heavily correlated with the class label. However, this color correlation is spurious as it has been reversed for test data. The correlations of labels with colors are +90% and +80% for training environments and -90% for the test environment.

Given the structure of data, if our learning algorithm only learns the spurious correlation (in this case color), the test accuracy will be 14%, while if our model learns the invariant features, we can achieve 75% test accuracy. The CMNIST setup (Arjovsky et al., 2019) we use has also been followed by Fishr (Rame et al., 2022) and V-Rex (Krueger et al., 2021). The only difference is that we replace the regularizer in IRM with our proposed regularizers in HGP and Hutchinson. In this setting, a simple fully connected network is used for training where we train the model for 501 steps and at step 190 the regularization parameters  $\alpha$  and  $\beta$  increase from 1 to 91257.18 suddenly (found by grid search for IRM). We repeated the experiments 10 times and reported the mean and average over these runs.

Table 2: Comparison of ERM, IRM, V-Rex, Fishr, and our proposed methods HGP and Hutchinson, on Colored MNIST experiment in introduced in IRM (Arjovsky et al., 2019) while the same hyperparameters have been used.

Method	Train acc.	Test acc.
ERM	86.4 ± 0.2	14.0 ± 0.7
IRM	71.0 ± 0.5	65.6 ± 1.8
V-REx	71.7 ± 1.5	67.2 ± 1.5
Fishr	71.0 ± 0.9	69.5 ± 1.0
HGP	71.0 ± 1.5	69.4 ± 1.3
Hutchinson	61.7 ± 1.9	<b>74.0 ± 1.2</b>

Table 2 compares the performance of HGP and Hutchinson against common DG baselines. As can be seen, while HGP achieves competitive performance against the state-of-the-art Fishr algorithm, Hutchinson outperforms all methods

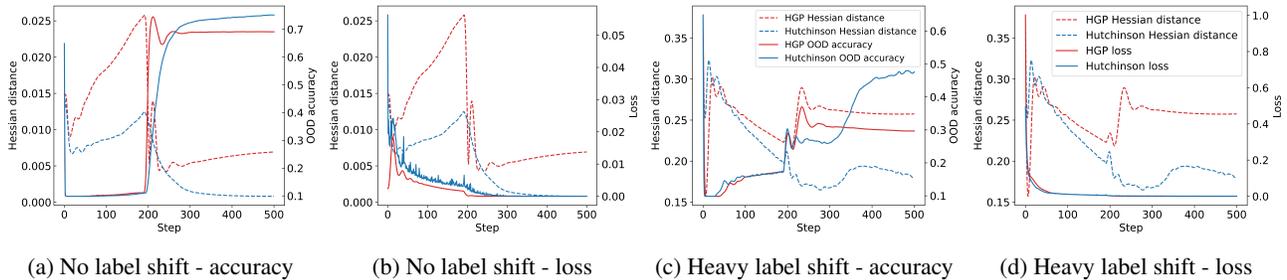


Figure 3: Correlation between Hessian distances ( $\|\mathbf{H}_{S_1} - \mathbf{H}_{S_2}\|_F$ ) and OOD accuracies/losses for HGP and Hutchinson regularization during the training for Colored MNIST. (a), (b): no label shift; (c), (d): heavy label shift.

by a large gap and obtains the near-optimal OOD accuracy of 74%, which is close to the maximum achievable accuracy of 75%.

### 5.3. OOD generalization under label shift

In order to explore the robustness of our proposed algorithm to label shift, we designed an imbalanced Colored MNIST dataset which is the same as the Colored MNIST dataset but we modified the dataset in a way that each of the two domains contain 95% of data points from one class and 5% from the other class. In fact, in the imbalanced Colored MNIST dataset, both correlation shift and label shift exist. We repeated the experiments in Table 2 on imbalanced Colored MNIST and reported the results in Table 3. Looking at Table 3, in the heavy label shift case, the Hutchinson achieves the highest OOD accuracy which outperforms the best-performing algorithm by a margin of 12%.

Table 3: Comparison of ERM, IRM, V-Rex, Fishr, and our proposed methods HGP and Hutchinson on Imbalanced Colored MNIST where each domain has 95% from one class and 5% from other class. Except for the imposed label shift, the setting is same as Colored MNIST experiment introduced in IRM (Arjovsky et al., 2019).

Method	Train acc.	Test acc.
ERM	$86.4 \pm 0.1$	$16.7 \pm 0.1$
IRM	$84.9 \pm 0.1$	$14.3 \pm 1.4$
V-REx	$83.3 \pm 0.2$	$35.1 \pm 1.2$
Fishr	$75.7 \pm 2.7$	$35.5 \pm 5.3$
HGP	$83.0 \pm 0.2$	$30.0 \pm 0.9$
Hutchinson	$79.4 \pm 0.3$	<b><math>47.7 \pm 1.4</math></b>

### 5.4. Training dynamics for Colored MNIST and Imbalanced Colored MNIST

Similar to earlier works (Arjovsky et al., 2019; Krueger et al., 2021; Rame et al., 2022), we study the training and OOD

generalization dynamics of HGP and Hutchinson on CMNIST and imbalanced CMNIST datasets. Figure 3a and Figure 3b show the Hessian distance  $\|\mathbf{H}_{S_1} - \mathbf{H}_{S_2}\|_F$  and OOD accuracy/loss of HGP and Hutchinson during the training. From Figure 3a when regularization parameters increase at step 190, the Hessian distance significantly drops and immediately OOD accuracy increases which suggests a strong negative correlation between the Hessian distance and OOD accuracy. This correlation is more evident for Hutchinson compared to HGP since after step 300, the Hessian distance keeps decreasing and OOD accuracy increases correspondingly. In contrast, for HGP, after step 300 when we are close to convergence, OOD accuracy becomes stable and even there is a slight increase in the Hessian distance. This is because close to convergence the gradients become zero and HGP loss cannot align Hessians as well as Hutchinson.

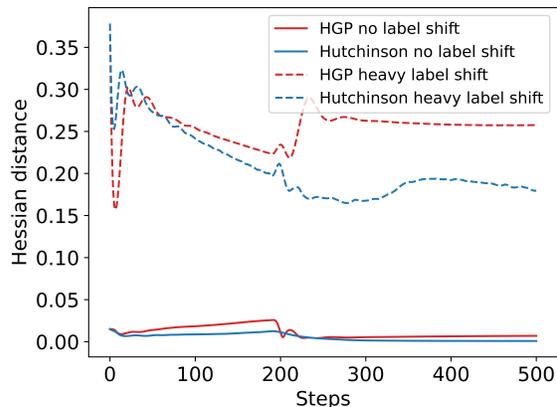


Figure 4: The Hessian distances  $\|\mathbf{H}_{S_1} - \mathbf{H}_{S_2}\|_F$  in Hutchinson and HGP regularization losses during the training steps for no label shift and heavy label shift in Colored MNIST experiment. Clearly, for heavy label shift, the Hessian distance is significantly more than no label shift case.

For the heavy label shift case, Figure 3c and 3d show that it is much more difficult to decrease the Hessian distance in

general. Besides, Hutchinson is more capable of reducing the Hessian distance compared to HGP and subsequently, it achieves better OOD accuracy.

Finally, Figure 4 compares Hessian distance while we use HGP and Hutchinson regularizers in no label shift and heavy label shift datasets. Clearly, for heavy label shift, the Hessian distance is significantly larger which shows that transferability between domains becomes more difficult. This experiment validates our argument that Hessian distance captures transferability between domains.

### 5.5. OOD generalization under diversity shift

In this section we use the Domainbed benchmark (Gulrajani and Lopez-Paz, 2020) to evaluate OOD accuracy of HGP and Hutchinson. We validate our proposed algorithms on the VLCS (Fang et al., 2013), PACS (Li et al., 2017), OfficeHome (Venkateswara et al., 2017), and DomainNet (Peng et al., 2019) datasets from the DomainBed benchmark. Domainbed provides a standard scheme to validate the performance of DG algorithms in a fair manner with different model selections. Table 4 provides the performance of HGP and Hutchinson methods against multiple benchmarks for leave-one-domain-out model selection reported from Domainbed. It shows that our Hutchinson method gives the state-of-the-art OOD performance in most cases. For HGP, our results are not as good as expected since the gradients diminish near convergence. The most competitive method to Hutchinson is CORAL, which we showed approximately aligns Hessians (by matching covariances) and gradients (by matching features). The results for other model selections are reported in the supplementary material.

### 5.6. Robustness to adversarial shift

Robustness of deep learning models to adversarial attacks is a notion of their generalization ability. To explore adversarial robustness of models trained using HGP and Hutchinson loss, we evaluate their performance against the Fast Gradient Sign Method (FGSM) attack (Goodfellow et al., 2014) and benchmark their performances against, ERM, IRM, V-Rex, and Fishr on CMNIST in Figure 5. We see that our Hutchinson method is also more transferable to adversarially perturbed domains.

### 5.7. Computational time comparison

For CMNIST experiments, we recorded the training time of our proposed algorithms (HGP and Hutchinson) and several popular domain generalization algorithms. The results of this experiment are presented in Table 5. We note that our Hessian alignment methods have more computation cost than existing algorithms, due to the Hessian estimation. This is the expense of better OOD generalization. However, the additional cost is controllable.

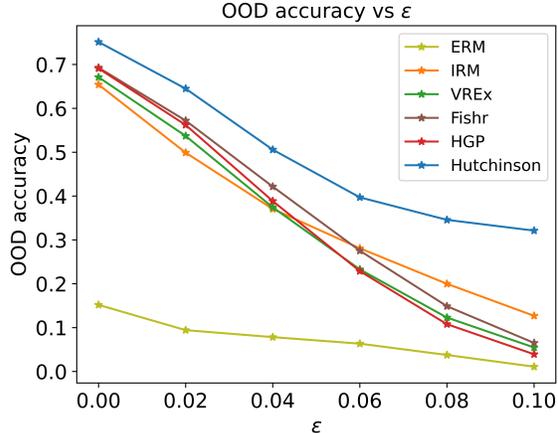


Figure 5: Comparison of OOD accuracy under adversarial attack.  $\epsilon$  denotes the perturbation amplitude.

### 5.8. The roles of Hessian and gradient alignment

Since we have implemented both Hessian and gradient alignment in our algorithms, which part is playing a main role? In this section, we conduct an ablation study to investigate the role of each regularization term in the loss functions of Hessian alignment. Recall that  $\alpha$  is the regularization weight of Hessian alignment and  $\beta$  is the weight of the gradient alignment. We set up an experiment with the PACS dataset and set the Sketch domain as test domain. We compare the Hutchinson algorithm in three different scenarios:

- Both:  $\alpha = 10^{-5}, \beta = 10^{-5}$ ;
- Only Hessian alignment:  $\alpha = 10^{-5}, \beta = 0$ ;
- Only gradient alignment:  $\alpha = 0, \beta = 10^{-5}$ .

The test domain accuracies for models trained with different hyperparameters are presented in Table 6. Clearly, with both Hessian and gradient alignment we achieve the best OOD generalization. Moreover, Hessian alignment plays a more important role than gradient alignment. Additional ablation studies for CMNIST experiment are presented in the appendix.

### 5.9. Implementation details

We provide the implementation details for all the experiments we conducted in the paper. For the transferability experiments in Figure 2, the experimental setup is exactly the same as the one used in Zhang et al. (2021). For regularization parameters, we used  $\alpha = 10^{-5}$  and  $\beta = 10^{-3}$  for HGP and Hutchinson methods. For the CMNIST experiment, the experimental setting is exactly the same as the IRM paper (Arjovsky et al., 2019). To set the regularization parameters, we also used the same scheme employed in

Table 4: DomainBed benchmark with leave-one-domain-out cross-validation model selection for CMNIST, VLCS, PACS, and OfficeHome datasets. We show the best and second best number with boldface and underline respectively. Due to time and computation limit, we did not finish the Fishr experiment on DomainNet.

Algorithm	VLCS	PACS	OfficeHome	DomainNet	Avg
ERM (Vapnik, 1999)	77.2	83.0	65.7	40.6	66.6
IRM (Arjovsky et al., 2019)	76.3	81.5	64.3	33.5	63.9
GroupDRO (Sagawa et al., 2020)	77.9	<u>83.5</u>	65.2	33.0	64.9
Mixup (Wang et al., 2020)	77.7	83.2	67.0	38.5	66.6
MLDG (Li et al., 2018a)	77.2	82.9	66.1	41.0	66.8
CORAL (Sun and Saenko, 2016)	<u>78.7</u>	82.6	<b>68.5</b>	<u>41.1</u>	<u>67.7</u>
MMD (Li et al., 2018b)	77.3	83.2	60.2	23.4	61.0
DANN (Ganin et al., 2016)	76.9	81.0	64.9	38.2	65.2
CDANN (Zhou et al., 2021)	77.5	78.8	64.3	38.0	64.6
MTL (Blanchard et al., 2021)	76.6	83.7	65.7	40.6	66.7
SagNet (Nam et al., 2020)	77.5	82.3	67.6	40.2	66.9
ARM (Zhang et al., 2020)	76.6	81.7	64.4	35.2	64.5
VREx (Krueger et al., 2021)	76.7	81.3	64.9	33.4	64.1
RSC (Huang et al., 2020)	77.5	82.6	65.8	38.9	66.2
Fishr (Rame et al., 2022)	78.2	<b>85.4</b>	<u>67.8</u>	-	-
HGP	76.7	82.2	67.5	<u>41.1</u>	66.9
Hutchinson	<b>79.3</b>	<u>84.8</u>	<b>68.5</b>	<b>41.4</b>	<b>68.5</b>

Table 5: Training time comparison on CMNIST between ERM, IRM, V-Rex, Fishr, and our proposed methods HGP and Hutchinson for 100 steps averaged over 10 runs.

Method	Train acc.
ERM	15.3 ± 0.1
CORAL	32.0 ± 3.2
V-REx	34.1 ± 2.3
IRM	27.0 ± 4.6
GroupDRO	32.8 ± 2.9
Fishr	25.3 ± 2.1
HGP	39.5 ± 4.4
Hutchinson	62.4 ± 1.5

Table 6: Ablation study of the Hutchinson method on different alignment regularization on the PACS dataset when the test domain is sketch.

Method	Test acc.
Hessian & gradient	81.4
Gradient only	77.0
Hessian only	79.4

the IRM paper (Arjovsky et al., 2019). The only difference here is that we have two regularization parameters so we set  $\alpha = \beta$ . For DomainBed experiments, the setup is the

same as in Gulrajani and Lopez-Paz (2020). For  $\alpha$  and  $\beta$ , we used ranges  $(-3, -5)$ ,  $(-1, -3)$  for HGP and  $(-1, -3)$ ,  $(-2, -4)$  for Hutchinson, respectively.

## 6. Conclusions

In this paper, we study the role Hessian alignment in domain generalization. Hessian alignment of the classifier head does not only improve the transferability to new domains, but also serves as an effective feature matching mechanism. In order to overcome the difficulty of heavy computation, we propose two estimation methods for Hessians, using Hessian-gradient product and Hutchinson’s diagonal estimation. These methods are tested to be effective in various scenarios, including spurious correlation (CMNIST), standard diversity shift benchmark (Domainbed), and adversarial shift. At the expense of slightly heavier computation, Hessian alignment performs competitively, often achieving the state-of-the-art. To our knowledge, our method provides the first Hessian alignment in domain generalization. In the future, it would be interesting to study and compare more efficient ways to align Hessians and gradients across domains.

## Acknowledgements

We want to thank the anonymous ICCV area chair and reviewers for constructive feedback. GZ would like to thank Han Zhao for valuable comments of the camera ready draft.

## References

- Arjovsky, M., Bottou, L., Gulrajani, I., and Lopez-Paz, D. (2019). Invariant risk minimization. *arXiv preprint arXiv:1907.02893*.
- Bekas, C., Kokiopoulou, E., and Saad, Y. (2007). An estimator for the diagonal of a matrix. *Applied numerical mathematics*, 57(11-12):1214–1229.
- Ben-David, S., Blitzer, J., Crammer, K., Kulesza, A., Pereira, F., and Vaughan, J. W. (2010). A theory of learning from different domains. *Machine learning*, 79(1):151–175.
- Blanchard, G., Deshmukh, A. A., Dogan, Ü., Lee, G., and Scott, C. (2021). Domain generalization by marginal transfer learning. *The Journal of Machine Learning Research*, 22(1):46–100.
- Chang, J.-R., Wu, M.-S., Yu, W.-H., Chen, C.-C., Yang, C.-K., Lin, Y.-Y., and Yeh, C.-Y. (2021). Stain Mix-Up: Unsupervised domain generalization for histopathology images. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 117–126. Springer.
- Fang, C., Xu, Y., and Rockmore, D. N. (2013). Unbiased metric learning: On the utilization of multiple datasets and web images for softening bias. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1657–1664.
- Ganin, Y., Ustinova, E., Ajakan, H., Germain, P., Larochelle, H., Laviolette, F., Marchand, M., and Lempitsky, V. (2016). Domain-adversarial training of neural networks. *The Journal of Machine Learning Research*, 17(1):2096–2030.
- Goodfellow, I. J., Shlens, J., and Szegedy, C. (2014). Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*.
- Gulrajani, I. and Lopez-Paz, D. (2020). In search of lost domain generalization. *arXiv preprint arXiv:2007.01434*.
- Hu, W., Niu, G., Sato, I., and Sugiyama, M. (2018). Does distributionally robust supervised learning give robust classifiers? In *International Conference on Machine Learning*, pages 2029–2037. PMLR.
- Huang, Z., Wang, H., Xing, E. P., and Huang, D. (2020). Self-challenging improves cross-domain generalization. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*, pages 124–140. Springer.
- Ilse, M., Tomczak, J. M., and Forré, P. (2021). Selecting data augmentation for simulating interventions. In *International Conference on Machine Learning*, pages 4555–4562. PMLR.
- Jaakkola, T. and Haussler, D. (1998). Exploiting generative models in discriminative classifiers. *Advances in Neural Information Processing Systems*, 11.
- Koltchinskii, V. (2010). Rademacher complexities and bounding the excess risk in active learning. *The Journal of Machine Learning Research*, 11:2457–2485.
- Koyama, M. and Yamaguchi, S. (2021). When is invariance useful in an out-of-distribution generalization problem?
- Krueger, D., Caballero, E., Jacobsen, J.-H., Zhang, A., Binas, J., Zhang, D., Le Priol, R., and Courville, A. (2021). Out-of-distribution generalization via risk extrapolation (REx). In *International Conference on Machine Learning*, pages 5815–5826. PMLR.
- Li, D., Yang, Y., Song, Y.-Z., and Hospedales, T. (2018a). Learning to generalize: Meta-learning for domain generalization. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32.
- Li, D., Yang, Y., Song, Y.-Z., and Hospedales, T. M. (2017). Deeper, broader and artier domain generalization. In *Proceedings of the IEEE international conference on computer vision*, pages 5542–5550.
- Li, H., Pan, S. J., Wang, S., and Kot, A. C. (2018b). Domain generalization with adversarial feature learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5400–5409.
- Mansilla, L., Echeveste, R., Milone, D. H., and Ferrante, E. (2021). Domain generalization via gradient surgery. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6630–6638.
- Muandet, K., Balduzzi, D., and Schölkopf, B. (2013). Domain generalization via invariant feature representation. In *International Conference on Machine Learning*, pages 10–18. PMLR.
- Nam, J., Cha, H., Ahn, S., Lee, J., and Shin, J. (2020). Learning from failure: De-biasing classifier from biased classifier. *Advances in Neural Information Processing Systems*, 33:20673–20684.
- Parascandolo, G., Neitz, A., Orvieto, A., Gresele, L., and Schölkopf, B. (2020). Learning explanations that are hard to vary. *arXiv preprint arXiv:2009.00329*.
- Peng, X., Bai, Q., Xia, X., Huang, Z., Saenko, K., and Wang, B. (2019). Moment matching for multi-source domain adaptation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1406–1415.
- Rame, A., Dancette, C., and Cord, M. (2022). Fishr: Invariant gradient variances for out-of-distribution generalization. In *International Conference on Machine Learning*, pages 18347–18377. PMLR.
- Sagawa, S., Koh, P. W., Hashimoto, T. B., and Liang, P. (2020). Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. In *International Conference on Learning Representations*.
- Shahtalebi, S., Gagnon-Audet, J.-C., Laleh, T., Faramarzi, M., Ahuja, K., and Rish, I. (2021). SAND-mask: An enhanced gradient masking strategy for the discovery of invariances in domain generalization. *arXiv preprint arXiv:2106.02266*.

- Shi, Y., Seely, J., Torr, P. H., Siddharth, N., Hannun, A., Usunier, N., and Synnaeve, G. (2021). Gradient matching for domain generalization. *arXiv preprint arXiv:2104.09937*.
- Skorski, M. (2021). Modern analysis of Hutchinson’s trace estimator. In *2021 55th Annual Conference on Information Sciences and Systems (CISS)*, pages 1–5. IEEE.
- Sun, B. and Saenko, K. (2016). Deep CORAL: Correlation alignment for deep domain adaptation. In *European conference on computer vision*, pages 443–450. Springer.
- Tzeng, E., Hoffman, J., Zhang, N., Saenko, K., and Darrell, T. (2014). Deep domain confusion: Maximizing for domain invariance. *arXiv preprint arXiv:1412.3474*.
- Vapnik, V. N. (1999). An overview of statistical learning theory. *IEEE Transactions on Neural Networks*, 10(5):988–999.
- Venkateswara, H., Eusebio, J., Chakraborty, S., and Panchanathan, S. (2017). Deep hashing network for unsupervised domain adaptation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5018–5027.
- Wang, J., Lan, C., Liu, C., Ouyang, Y., Qin, T., Lu, W., Chen, Y., Zeng, W., and Yu, P. (2022). Generalizing to unseen domains: A survey on domain generalization. *IEEE Transactions on Knowledge and Data Engineering*.
- Wang, Y., Li, H., and Kot, A. C. (2020). Heterogeneous domain generalization via domain mixup. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3622–3626. IEEE.
- Ye, N., Li, K., Hong, L., Bai, H., Chen, Y., Zhou, F., and Li, Z. (2021). OoD-Bench: Benchmarking and understanding out-of-distribution generalization datasets and algorithms. *arXiv preprint arXiv:2106.03721*, 1(3):5.
- Zhang, G., Zhao, H., Yu, Y., and Poupart, P. (2021). Quantifying and improving transferability in domain generalization. *Advances in Neural Information Processing Systems*, 34:10957–10970.
- Zhang, H., Cisse, M., Dauphin, Y. N., and Lopez-Paz, D. (2018). mixup: Beyond empirical risk minimization. In *International Conference on Learning Representations*.
- Zhang, M., Marklund, H., Gupta, A., Levine, S., and Finn, C. (2020). Adaptive risk minimization: A meta-learning approach for tackling group shift. *arXiv preprint arXiv:2007.02931*, 8:9.
- Zhao, H., Des Combes, R. T., Zhang, K., and Gordon, G. (2019). On learning invariant representations for domain adaptation. In *International conference on machine learning*, pages 7523–7532. PMLR.
- Zhao, H., Zhang, S., Wu, G., Moura, J. M., Costeira, J. P., and Gordon, G. J. (2018). Adversarial multiple source domain adaptation. *Advances in neural information processing systems*, 31.
- Zhou, K., Liu, Z., Qiao, Y., Xiang, T., and Loy, C. C. (2021). Domain generalization in vision: A survey. *arXiv preprint arXiv:2103.02503*.

## A. Proof of Proposition 4

*Proof.* To formulate the classifier head of the neural network, let  $z_i$  be the  $i$ -th component of the feature vector before the classifier layer and the classifier's parameter  $\theta$  is decomposed to  $w_{k,i}$ , the element in row  $k$  and column  $i$  of the classifier weight matrix, and  $b_k$ , the bias term for the  $k$ -th output. We define  $a_k$  as

$$a_k = \sum_i^c w_{k,i} z_i + b_k, \quad (24)$$

where  $c$  is the number of classes. Given  $a_k$ , if we assume the classifier activation function  $\sigma(\cdot)$  to be softmax, the classifier output for the  $k$ -th neuron can be written as

$$\hat{y}_k = \sigma(a_k) = \frac{e^{a_k}}{\sum_{j=1}^c e^{a_j}}, \quad (25)$$

Now, if we denote  $(\mathbf{x}, \mathbf{y})$  as the sample and  $\hat{\mathbf{y}}$  as the associated output, and assume the loss function be cross-entropy

$$\ell(\mathbf{x}, \mathbf{y}; \theta) = - \sum_{l=1}^c y_l \log \hat{y}_l, \quad (26)$$

Having the loss function, we proceed to calculate the gradients and hessian of loss with respect to classifier parameters  $w_{p,q}$  and  $b_p$  is

$$\frac{\partial \ell(\hat{\mathbf{y}}, \mathbf{y}; \theta)}{\partial w_{p,q}} = \sum_{l=1}^c \frac{\partial \ell}{\partial \hat{y}_l} \sum_{k=1}^c \frac{\partial \hat{y}_l}{\partial a_k} \frac{\partial a_k}{\partial w_{p,q}}. \quad (27)$$

$$\frac{\partial \ell(\hat{\mathbf{y}}, \mathbf{y}; \theta)}{\partial b_u} = \sum_{l=1}^c \frac{\partial \ell}{\partial \hat{y}_l} \sum_{k=1}^c \frac{\partial \hat{y}_l}{\partial a_k} \frac{\partial a_k}{\partial b_u}. \quad (28)$$

Given that the activation function is a softmax function  $\hat{y}_l = \sigma(a_l) = \frac{e^{a_l}}{\sum_j e^{a_j}}$ , the  $\frac{\partial \sigma(a_l)}{\partial a_k}$  can be calculated as:

$$\frac{\partial \sigma(a_l)}{\partial a_k} = \sigma(a_l) \delta_{l,k} - \sigma(a_l) \sigma(a_k). \quad (29)$$

Now, using Eq. 29, we can rewrite Eqs. 27 and 28 as follows:

$$\frac{\partial \ell}{\partial w_{p,q}} = (\hat{y}_p - y_p) z_q, \quad \frac{\partial \ell}{\partial b_u} = (\hat{y}_u - y_u). \quad (30)$$

For the Hessian, we only calculate the elements of matrix that are only related to the classifier layer. More precisely, we calculate  $\frac{\partial^2 \ell}{\partial w_{u,v} \partial w_{p,q}}$ ,  $\frac{\partial^2 \ell}{\partial w_{p,q} \partial b_u}$ , and  $\frac{\partial^2 \ell}{\partial b_u \partial b_v}$ . For the  $\frac{\partial^2 \ell}{\partial w_{u,v} \partial w_{p,q}}$  we can write

$$\frac{\partial^2 \ell}{\partial w_{u,v} \partial w_{p,q}} = \frac{\partial}{\partial w_{u,v}} ((\hat{y}_p - y_p) z_q). \quad (31)$$

To calculate the above expression, we need  $\frac{\partial \hat{y}_p}{\partial w_{u,v}}$ :

$$\frac{\partial \hat{y}_p}{\partial w_{u,v}} = \sum_k \frac{\partial \hat{y}_p}{\partial a_k} \frac{\partial a_k}{\partial w_{u,v}} = \sum_k \frac{\partial \sigma(a_p)}{\partial a_k} \sum_i \delta_{k,u} \delta_{i,v} z_i = \sum_k \sigma(a_p) (\delta_{p,k} - \sigma(a_k)) \delta_{k,u} z_v = \hat{y}_p z_v (\delta_{p,u} - \hat{y}_u) \quad (32)$$

Now, incorporating eq. 32 into eq. 31, the elements of Hessian matrix in classifier layer i.e.,  $\frac{\partial^2 \ell}{\partial w_{u,v} \partial w_{p,q}}$ , we have:

$$\frac{\partial^2 \ell}{\partial w_{u,v} \partial w_{p,q}} = z_q z_v \hat{y}_p (\delta_{p,u} - \hat{y}_u). \quad (33)$$

For  $\frac{\partial^2 \ell}{\partial w_{p,q} \partial b_u}$  we can write

$$\frac{\partial^2 \ell}{\partial w_{p,q} \partial b_u} = \frac{\partial}{\partial b_u} ((\hat{y}_p - y_p) z_q) = \frac{\partial(\hat{y}_p - y_p)}{\partial b_u} z_q = \sum_k \frac{\partial \hat{y}_p}{\partial a_k} \frac{\partial a_k}{\partial b_u} z_q = z_q \hat{y}_p (\delta_{p,u} - \hat{y}_u). \quad (34)$$

Eventually, for  $\frac{\partial^2 \mathcal{L}}{\partial b_u \partial b_v}$  we have:

$$\frac{\partial^2 \mathcal{L}}{\partial b_u \partial b_v} = \frac{\partial}{\partial b_v} (\hat{y}_u - y_u) = \sum_k \frac{\partial \hat{y}_u}{\partial a_k} \frac{\partial a_k}{\partial b_v} = \hat{y}_u (\delta_{u,v} - \hat{y}_v). \quad (35)$$

□

## B. Alignment attributes in Hessian and gradient in regression task with mean square error loss and general activation function

In this section, we extend our analysis to regression tasks for real numbers. We show that the classifier head's gradient and Hessian yield similar information of the features. To adapt our framework to regression, we replace the meaning of  $\sigma$  from softmax to an arbitrary uni-variate activation function.

**Proposition 5 (Alignment attributes in Hessian and gradient for mean square error loss and general activation function).**

Let  $\hat{y}$  and  $y$  be the network prediction and true target associated with the output neuron of a single output network,  $\sigma(\cdot)$  be the activation function,  $z_i$  be the  $i$ -th feature value before the last layer (regression layer). Suppose the last layer's parameter  $\theta$  is decomposed to  $w_i$ , the  $i$ -th element of the weight vector, and  $b$ , the bias term. Matching the gradients and Hessians with respect to the last layer across the domain aligns the following attributes

$$\frac{\partial \ell}{\partial b} = (\hat{y} - y) \sigma'(a), \quad (36)$$

$$\frac{\partial \ell}{\partial w_i} = (\hat{y} - y) \sigma'(a) z_i, \quad (37)$$

$$\frac{\partial^2 \ell}{\partial b^2} = \sigma'(a)^2 + (\hat{y} - y) \sigma''(a), \quad (38)$$

$$\frac{\partial^2 \ell}{\partial w_i \partial b} = \sigma'(a)^2 z_i + (\hat{y} - y) \sigma''(a) z_i, \quad (39)$$

$$\frac{\partial^2 \ell}{\partial w_i \partial w_k} = \sigma'(a)^2 z_i z_k + (\hat{y} - y) \sigma''(a) z_i z_k, \quad (40)$$

*Proof.* To formulate the last layer of the neural network we define  $a$  as

$$a = \sum_i w_i z_i + b. \quad (41)$$

Given  $a$ , if we assume the last layer activation function is  $\sigma(\cdot)$ , the single output can be written as

$$\hat{y} = \sigma(a). \quad (42)$$

Now, if we denote the input data as  $(\mathbf{x}, y)$  and associated output  $\hat{y}$ , and assume the loss function be

$$\ell(\mathbf{x}, y; \theta) = \frac{1}{2} (\hat{y} - y)^2. \quad (43)$$

Having the loss function, we proceed to calculate the gradients and hessian of loss with respect to last layer parameters  $w_i$  and  $b$ . For the gradients, we write

$$\frac{\partial \ell}{\partial w_i} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a} \frac{\partial a}{\partial w_i} = (\hat{y} - y) \sigma'(a) z_i, \quad (44)$$

$$\frac{\partial \ell}{\partial b} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a} \frac{\partial a}{\partial b} = (\hat{y} - y) \sigma'(a). \quad (45)$$

For the Hessian matrix, we only calculate the elements of the matrix that are only related to the last layer. More precisely, we calculate  $\frac{\partial^2 \mathcal{L}}{\partial w_i \partial w_k}$ ,  $\frac{\partial^2 \mathcal{L}}{\partial w_i \partial b}$ , and  $\frac{\partial^2 \mathcal{L}}{\partial b^2}$ . For the  $\frac{\partial^2 \mathcal{L}}{\partial w_i \partial w_k}$  we can write

$$\frac{\partial^2 \ell}{\partial w_i \partial w_k} = \frac{\partial}{\partial w_k} ((\hat{y} - y) \cdot \sigma'(a) z_i) = \frac{\partial(\hat{y} - y)}{\partial w_k} \cdot \sigma'(a) z_i + (\hat{y} - y) \frac{\partial}{\partial w_k} \sigma'(a) z_i \quad (46)$$

To calculate the above expression, we need  $\frac{\partial \hat{y}}{\partial w_k}$  and  $\frac{\partial}{\partial w_k} \sigma'(a)$ .

$$\frac{\partial \hat{y}}{\partial w_k} = \frac{\partial \hat{y}}{\partial a} \frac{\partial a}{\partial w_k} = \sigma'(a) z_k, \quad (47)$$

$$\frac{\partial}{\partial w_k} \sigma'(a) = \frac{\partial}{\partial a} \sigma'(a) \frac{\partial a}{\partial w_k} = \sigma''(a) z_k. \quad (48)$$

Now, incorporating eq. 48 into 46 the elements of the last-layer Hessian matrix become:

$$\frac{\partial^2 \ell}{\partial w_i \partial w_k} = \sigma'(a)^2 z_i z_k + (\hat{y} - y) \sigma''(a) z_i z_k. \quad (49)$$

For  $\frac{\partial^2 \mathcal{L}}{\partial w_i \partial b}$  we can write

$$\frac{\partial^2 \ell}{\partial w_i \partial b} = \frac{\partial}{\partial b} ((\hat{y} - y) \cdot \sigma'(a) z_i) = \frac{\partial(\hat{y} - y)}{\partial b} \cdot \sigma'(a) z_i + (\hat{y} - y) \frac{\partial}{\partial b} \sigma'(a) z_i \quad (50)$$

$$= \frac{\partial \hat{y}}{\partial a} \frac{\partial a}{\partial b} \cdot \sigma'(a) z_i + (\hat{y} - y) \frac{\partial}{\partial a} \sigma'(a) \frac{\partial a}{\partial b} z_i = \sigma'(a)^2 z_i + (\hat{y} - y) \sigma''(a) z_i. \quad (51)$$

Eventually,  $\frac{\partial^2 \mathcal{L}}{\partial b^2}$  is calculated as:

$$\frac{\partial^2 \ell}{\partial b^2} = \frac{\partial}{\partial b} (\hat{y} - y) \cdot \sigma'(a) = \frac{\partial(\hat{y} - y)}{\partial b} \cdot \sigma'(a) + (\hat{y} - y) \frac{\partial}{\partial b} \sigma'(a) \quad (52)$$

$$= \frac{\partial \hat{y}}{\partial a} \frac{\partial a}{\partial b} \cdot \sigma'(a) + (\hat{y} - y) \frac{\partial}{\partial a} \sigma'(a) \frac{\partial a}{\partial b} = \sigma'(a)^2 + (\hat{y} - y) \sigma''(a). \quad (53)$$

□

Eqs. 45, 49, 51, 53 show that matching Hessians and gradients with respect to the last layer parameters will match neural network outputs, last layer input features and covariance between output features across domains. This supports the idea of matching gradients and Hessians during training. The above result can be extended to multi-dimensional outputs if the activation is element-wise.

## C. Ablation Study

We also repeat the Colored MNIST and imbalanced Colored MNIST experiments in scenarios where one of  $\alpha$  and  $\beta$  is non-zero, in Table 7 and Table 8. Recall that  $\alpha$  controls the Hessian alignment and  $\beta$  controls the gradient alignment. If not mentioned, the values for  $\alpha$  and/or  $\beta$  are non-zero and they are chosen exactly as the IRM paper (Arjovsky et al., 2019). According to Table 7, for both HGP and Hutchinson methods, gradient alignment seems to contribute more to OOD generalization on CMNIST. This might be due to the heavy correlation shift and we have to align the local minima first. For Hutchinson, when  $\beta = 0$ , the OOD performance drops which we believe is because the value that has been chosen for  $\alpha$  is optimized for the IRM loss. In other words, if we optimize  $\alpha$  for aligning the diagonal part of Hessian, it can contribute to the OOD generalization. For imbalanced Colored MNIST, the same trend for the role of  $\alpha$  and  $\beta$  can be observed. Overall, the key observation is that both aligning gradients and diagonal parts of Hessians contribute to the OOD generalization.

Table 7: Comparison of ERM, IRM, V-Rex, Fishr, and our proposed methods HGP and Hutchinson with ablation study on  $\alpha$  and  $\beta$  on Colored MNIST. The setting is same as the Colored MNIST experiment introduced in IRM (Arjovsky et al., 2019).

Method	Train acc.	Test acc.
ERM	86.4 $\pm$ 0.2	14.0 $\pm$ 0.7
IRM	71.0 $\pm$ 0.5	65.6 $\pm$ 1.8
V-REx	71.7 $\pm$ 1.5	67.2 $\pm$ 1.5
Fishr	71.0 $\pm$ 0.9	69.5 $\pm$ 1.0
HGP	71.0 $\pm$ 1.5	69.4 $\pm$ 1.3
HGP ( $\alpha = 0$ )	70.6 $\pm$ 1.8	69.3 $\pm$ 1.2
HGP ( $\beta = 0$ )	78.9 $\pm$ 0.3	53.3 $\pm$ 1.7
Hutchinson	61.7 $\pm$ 1.9	<b>74.0 <math>\pm</math> 1.2</b>
Hutchinson ( $\alpha = 0$ )	70.6 $\pm$ 1.8	69.3 $\pm$ 1.2
Hutchinson ( $\beta = 0$ )	84.9 $\pm$ 0.1	9.8 $\pm$ 0.2

Table 8: Comparison of ERM, IRM, V-Rex, Fishr, and our proposed methods HGP and Hutchinson with ablation study on  $\alpha$  and  $\beta$  on Imbalanced Colored MNIST where each domain has 95% from one class and 5% from other class. Except for the imposed label shift, the setting is same as the Colored MNIST experiment introduced in IRM (Arjovsky et al., 2019).

Method	Train acc.	Test acc.
ERM	86.4 $\pm$ 0.1	16.7 $\pm$ 0.1
IRM	84.9 $\pm$ 0.1	14.3 $\pm$ 1.4
V-REx	83.3 $\pm$ 0.2	35.1 $\pm$ 1.2
Fishr	75.7 $\pm$ 2.7	35.5 $\pm$ 5.3
HGP	83.0 $\pm$ 0.2	30.0 $\pm$ 0.9
HGP ( $\alpha = 0$ )	83.2 $\pm$ 0.4	29.7 $\pm$ 1.7
HGP ( $\beta = 0$ )	84.3 $\pm$ 0.1	20.4 $\pm$ 1.0
Hutchinson	79.4 $\pm$ 0.3	<b>47.7 <math>\pm</math> 1.4</b>
Hutchinson ( $\alpha = 0$ )	83.2 $\pm$ 0.4	29.7 $\pm$ 1.7
Hutchinson ( $\beta = 0$ )	84.9 $\pm$ 0.1	9.8 $\pm$ 0.2

## D. Domainbed Results for Other Model Selection Methods

In this section, we provide the Domainbed results for the two other model selection methods, i.e., the training-domain validation set and test-domain validation set (oracle). First note that the oracle model selection is not a valid benchmarking scheme and not applicable in practice as it uses the target domain data for selecting the hyperparameters. In fact, in this scenario, algorithms with more hyperparameters and training tricks (like warmup, exponential moving average and etc) can obtain better performance since they have more freedom to tune the model on test data. Considering this, we should not rely on the oracle model selection technique to compare domain generalization algorithms. The other model selection technique is the training-domain validation set where the validation sets of all training domain are concatenated together and select the hyperparameters that maximize the accuracy on the entire validation set.

As can be seen in Table 9 and Table 10, although Hessian alignment methods are not the best, their performance across all datasets is still competitive for other model selections. As also shown in Gulrajani and Lopez-Paz (2020), different model selections could result in different rankings of the algorithms. We find that training-domain model selection in general gives better results for most baseline algorithms, but the performance of the Hutchinson method slightly degrades. We defer the study of comparing model selections to future work.

Table 9: DomainBed benchmark with *training-domain validation set model selection* method for CMNIST, VLCS, PACS, and OfficeHome datasets. We show the best and second best number with boldface and underline respectively.

Algorithm	VLCS	PACS	OfficeHome	DomainNet	Avg
ERM	77.5	85.5	66.5	40.9	67.6
IRM	78.5	83.5	64.3	33.9	65.1
GroupDRO	76.7	84.4	66.0	33.3	65.1
Mixup	77.4	84.6	68.1	39.2	67.3
MLDG	77.2	84.9	66.8	41.2	67.5
CORAL	<b>78.8</b>	<u>86.2</u>	<b>68.7</b>	41.5	<b>68.8</b>
MMD	77.5	84.6	66.3	23.4	63.0
DANN	<u>78.6</u>	83.6	65.9	38.3	66.6
CDANN	77.5	82.6	65.8	38.3	66.1
MTL	77.2	84.6	66.4	40.6	67.2
SagNet	77.8	<b>86.3</b>	68.1	40.3	68.1
ARM	77.6	85.1	64.8	35.5	65.8
VREx	78.3	84.9	66.4	33.6	65.8
RSC	77.1	85.2	65.5	38.9	66.7
AND-mask	78.1	84.4	65.5	37.2	66.3
SAND-mask	77.4	84.6	65.8	32.1	65.0
Fish	77.8	85.5	<u>68.6</u>	<b>42.7</b>	<u>68.7</u>
Fishr	77.8	85.8	67.8	<u>41.7</u>	68.3
HGP	77.6	84.7	68.2	41.1	67.9
Hutchinson	76.8	83.9	68.2	41.6	67.6

Table 10: DomainBed benchmark with *test-domain validation set (oracle)* model selection method for CMNIST, VLCS, PACS, and OfficeHome datasets. We show the best and second best number with boldface and underline respectively.

Algorithm	VLCS	PACS	OfficeHome	DomainNet	Avg
ERM	77.6	86.7	66.4	41.3	68.0
IRM	76.9	84.5	63.0	28.0	63.1
GroupDRO	77.4	<u>87.1</u>	66.2	33.4	66.0
Mixup	78.1	86.8	68.0	39.6	68.1
MLDG	77.5	86.8	66.6	41.6	68.1
CORAL	77.7	<u>87.1</u>	<b>68.4</b>	41.8	<u>68.8</u>
MMD	77.9	<b>87.2</b>	66.2	23.5	63.7
DANN	<u>79.7</u>	85.2	65.3	38.3	67.1
CDANN	<b>79.9</b>	85.8	65.3	38.5	67.4
MTL	77.7	86.7	66.5	40.8	67.9
SagNet	77.6	86.4	67.5	40.8	68.1
ARM	77.8	85.8	64.8	36.0	66.1
VREx	78.1	<b>87.2</b>	65.7	30.1	65.3
RSC	77.8	86.2	66.5	38.9	67.4
AND-mask	76.4	86.4	66.1	37.9	66.7
SAND-mask	76.2	85.9	65.9	32.2	65.1
Fish	77.8	85.8	66.0	<u>42.7</u>	68.1
Fishr	78.2	86.9	<u>68.2</u>	<b>43.4</b>	<b>69.2</b>
HGP	77.3	86.5	67.4	41.2	68.1
Hutchinson	77.9	86.3	<b>68.4</b>	41.9	68.6