



# Interventional Few-Shot Learning

Zhongqi Yue<sup>1,3</sup>, Hanwang Zhang<sup>1</sup>, Qianru Sun<sup>2</sup>, Xian-Sheng Hua<sup>3</sup>

<sup>1</sup>Nanyang Technological University, <sup>2</sup>Singapore Management University,

<sup>3</sup>Damo Academy, Alibaba Group

yuez0003@ntu.edu.sg, hanwangzhang@ntu.edu.sg,

qianrusun@smu.edu.sg, xiansheng.hxs@alibaba-inc.com

## Abstract

We uncover an ever-overlooked deficiency in the prevailing Few-Shot Learning (FSL) methods: the pre-trained knowledge is indeed a confounder that limits the performance. This finding is rooted from our causal assumption: a Structural Causal Model (SCM) for the causalities among the pre-trained knowledge, sample features, and labels. Thanks to it, we propose a novel FSL paradigm: Interventional Few-Shot Learning (IFSL). Specifically, we develop three effective IFSL algorithmic implementations based on the backdoor adjustment, which is essentially a causal intervention towards the SCM of many-shot learning: the upper-bound of FSL in a causal view. It is worth noting that the contribution of IFSL is orthogonal to existing fine-tuning and meta-learning based FSL methods, hence IFSL can improve all of them, achieving a new 1-/5-shot state-of-the-art on *miniImageNet*, *tieredImageNet*, and cross-domain CUB. Code is released at <https://github.com/yue-zhongqi/ifs1>.

## 1 Introduction

Few-Shot Learning (FSL) — the task of training a model using very few samples — is nothing short of a panacea for any scenario that requires fast model adaptation to new tasks [69], such as minimizing the need for expensive trials in reinforcement learning [31] and saving computation resource for light-weight neural networks [28, 26]. Although we knew that, more than a decade ago, the crux of FSL is to imitate the human ability of transferring prior knowledge to new tasks [19], not until the recent advances in pre-training techniques, had we yet reached a consensus on “what & how to transfer”: a powerful neural network  $\Omega$  pre-trained on a large dataset  $\mathcal{D}$ . In fact, the prior knowledge learned from pre-training prospers today’s deep learning era, *e.g.*,  $\mathcal{D}$  = ImageNet,  $\Omega$  = ResNet in visual recognition [25, 24];  $\mathcal{D}$  = Wikipedia,  $\Omega$  = BERT in natural language processing [65, 17].

In the context of pre-trained knowledge, we denote the original FSL training set as *support* set  $\mathcal{S}$  and the test set as *query* set  $\mathcal{Q}$ , where the classes in  $(\mathcal{S}, \mathcal{Q})$  are unseen (or new) in  $\mathcal{D}$ . Then, we can use  $\Omega$  as a backbone (fixed or partially trainable) for extracting sample representations  $\mathbf{x}$ , and thus FSL can be achieved simply by *fine-tuning* the target model on  $\mathcal{S}$  and test it on  $\mathcal{Q}$  [13, 18]. However, the fine-tuning only exploits the  $\mathcal{D}$ ’s knowledge on “what to transfer”, but neglects “how to transfer”. Fortunately, the latter can be addressed by applying a post-pre-training and pre-fine-tuning strategy: *meta-learning* [55]. Different from fine-tuning

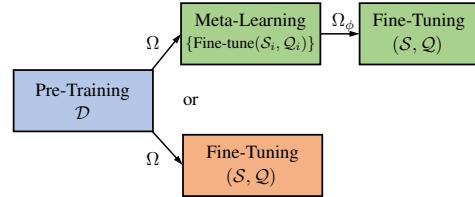


Figure 1: The relationships among different FSL paradigms (color green and orange). Our goal is to remove the deficiency introduced by Pre-Training.

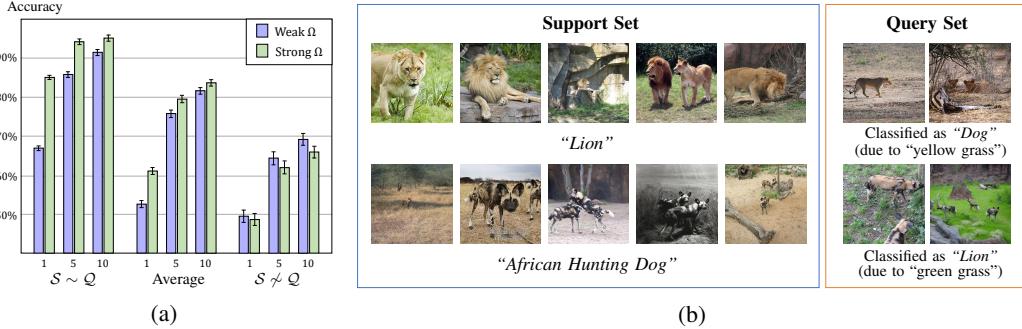


Figure 2: Quantitative and qualitative evidences of pre-trained knowledge misleading the fine-tune FSL paradigm. (a) *miniImageNet* fine-tuning accuracy on 1-/5-/10-shot FSL using weak and strong backbones: ResNet-10 and WRN-28-10.  $\mathcal{S} \sim \mathcal{Q}$  ( $\mathcal{S} \neq \mathcal{Q}$ ) denotes the pre-trained classifier scores of the query is similar (or dissimilar) to that of the support set. ‘‘Average’’ is the mean of both. The dissimilarity is measured using query hardness defined in Section 5.1. (b) An example of 5-shot  $\mathcal{S} \neq \mathcal{Q}$ .

whose goal is the ‘‘model’’ trained on  $\mathcal{S}$  and tested on  $\mathcal{Q}$ , meta-learning aims to learn the ‘‘meta-model’’ — a learning behavior — trained on many learning episodes  $\{(\mathcal{S}_i, \mathcal{Q}_i)\}$  sampled from  $\mathcal{D}$  and tested on the target task  $(\mathcal{S}, \mathcal{Q})$ . In particular, the behavior can be parametrized by  $\phi$  using model parameter generator [49, 21] or initialization [20]. After meta-learning, we denote  $\Omega_\phi$  as the new model starting point for the subsequent fine-tuning on target task  $(\mathcal{S}, \mathcal{Q})$ . Figure 1 illustrates the relationships among the above discussed FSL paradigms.

It is arguably a common sense that the stronger the pre-trained  $\Omega$  is, the better the downstream model will be. However, we surprisingly find that this may not be always the case in FSL. As shown in Figure 2(a), we can see a paradox: though stronger  $\Omega$  improves the performance on average, it indeed degrades that of samples in  $\mathcal{Q}$  dissimilar to  $\mathcal{S}$ . To illustrate the ‘‘dissimilar’’, we show a 5-shot learning example in Figure 2(b), where the prior knowledge on ‘‘green grass’’ and ‘‘yellow grass’’ is misleading. For example, the ‘‘Lion’’ samples in  $\mathcal{Q}$  have ‘‘yellow grass’’, hence they are misclassified as ‘‘Dog’’ whose  $\mathcal{S}$  has major ‘‘yellow grass’’. If we use stronger  $\Omega$ , the seen old knowledge (‘‘grass’’ & ‘‘color’’) will be more robust than the unseen new knowledge (‘‘Lion’’ & ‘‘Dog’’), and thus the old becomes even more misleading. We believe that such a paradox reveals an unknown systematic deficiency in FSL, which has been however hidden for years by our gold-standard ‘‘fair’’ accuracy, averaged over all the random  $(\mathcal{S}, \mathcal{Q})$  test trials, regardless of the similarity between  $\mathcal{S}$  and  $\mathcal{Q}$  (cf. Figure 2(a)). Though Figure 2 only illustrates the fine-tune FSL paradigm, the deficiency is expected in the meta-learning paradigm, as fine-tune is also used in each meta-train episode (Figure 1). We will analyze them thoroughly in Section 5.

In this paper, we first point out that the cause of the deficiency: pre-training can do evil in FSL, and then propose a novel FSL paradigm: Interventional Few-Shot Learning (IFSL), to counter the evil. Our theory is based on the assumption of the *causalities* among the pre-trained knowledge, few-shot samples, and class labels. Specifically, our contributions are summarized as follows.

- We begin with a Structural Causal Model (SCM) assumption in Section 2.2, which shows that the pre-trained knowledge is essentially a *confounder* that causes spurious correlations between the sample features and class labels in support set. As an intuitive example in Figure 2(b), even though the ‘‘grass’’ feature is not the cause of the ‘‘Lion’’ label, the prior knowledge on ‘‘grass’’ still confounds the classifier to learn a correlation between them.
- In Section 2.3, we illustrate a causal justification of why the proposed IFSL fundamentally works better: it is essentially a causal approximation to many-shot learning. This motivates us to develop three effective implementations of IFSL using the backdoor adjustment [46] in Section 3.
- Thanks to the causal intervention, IFSL is naturally orthogonal to the downstream fine-tuning and meta-learning based FSL methods [20, 66, 29]. In Section 5.2, IFSL improves all baselines by a considerable margin, achieving new 1-/5-shot state-of-the-arts: 73.51%/83.21% on *miniImageNet* [66], 83.07%/88.69% on *tieredImageNet* [52], and 50.71%/64.43% on cross-domain CUB [70].
- We further diagnose the detailed performances of FSL methods across different similarities between  $\mathcal{S}$  and  $\mathcal{Q}$ . We find that IFSL outperforms all baselines in every inch.

## 2 Problem Formulations

### 2.1 Few-Shot Learning

We are interested in a prototypical FSL: train a  $K$ -way classifier on an  $N$ -shot support set  $\mathcal{S}$ , where  $N$  is a small number of training samples per class (e.g.,  $N=1$  or  $5$ ); then test the classifier on a query set  $\mathcal{Q}$ . As illustrated in Figure 1, we have the following two paradigms to train the classifier  $P(y|\mathbf{x}; \theta)$ , predicting the class  $y \in \{1, \dots, K\}$  of a sample  $\mathbf{x}$ :

**Fine-Tuning.** We consider the prior knowledge as the sample feature representation  $\mathbf{x}$ , encoded by the pre-trained network  $\Omega$  on dataset  $\mathcal{D}$ . In particular, we refer  $\mathbf{x}$  to the output of the frozen sub-part of  $\Omega$  and the rest trainable sub-part of  $\Omega$  (if any) can be absorbed into  $\theta$ . We train the classifier  $P(y|\mathbf{x}; \theta)$  on the support set  $\mathcal{S}$ , and then evaluate it on the query set  $\mathcal{Q}$  in a standard supervised way.

**Meta-Learning.** Yet,  $\Omega$  only carries prior knowledge in a way of “representation”. If the dataset  $\mathcal{D}$  can be re-organized as the training episodes  $\{(\mathcal{S}_i, \mathcal{Q}_i)\}$ , each of which can be treated as a “sandbox” that has the same  $N$ -shot- $K$ -way setting as the target  $(\mathcal{S}, \mathcal{Q})$ . Then, we can model the “learning behavior” from  $\mathcal{D}$  parameterized as  $\phi$ , which can be learned by the above fine-tuning paradigm for each  $(\mathcal{S}_i, \mathcal{Q}_i)$ . Formally, we denote  $P_\phi(y|\mathbf{x}; \theta)$  as the enhanced classifier equipped with the learned behavior. For example,  $\phi$  can be the classifier weight generator [21], distance kernel function in  $k$ -NN [66], or even  $\theta$ 's initialization [20]. Considering  $L_\phi(\mathcal{S}_i, \mathcal{Q}_i; \theta)$  as the loss function of  $P_\phi(y|\mathbf{x}; \theta)$  trained on  $\mathcal{S}_i$  and tested on  $\mathcal{Q}_i$ , we can have  $\phi \leftarrow \arg \min_{(\phi, \theta)} \mathbb{E}_i [L_\phi(\mathcal{S}_i, \mathcal{Q}_i; \theta)]$ , and then we fix the optimized  $\phi$  and fine-tune for  $\theta$  on  $\mathcal{S}$  and test on  $\mathcal{Q}$ . Please refer to Appendix 5 for the details of various fine-tuning and meta-learning settings.

### 2.2 Structural Causal Model

From the above discussion, we can see that  $(\phi, \theta)$  in meta-learning and  $\theta$  in fine-tuning are both dependent on the pre-training. Such “dependency” can be formalized with a Structural Causal Model (SCM) [46] proposed in Figure 3(a), where the nodes denote the abstract data variables and the directed edges denote the (functional) causality, e.g.,  $X \rightarrow Y$  denotes that  $X$  is the cause and  $Y$  is the effect. Now we introduce the graph and the rationale behind its construction at a high-level. Please see Section 3 for the detailed functional implementations.

$D \rightarrow X$ . We denote  $X$  as the feature representation and  $D$  as the pre-trained knowledge, e.g., the dataset  $\mathcal{D}$  and its induced model  $\Omega$ . This link assumes that the feature  $X$  is extracted by using  $\Omega$ .

$D \rightarrow C \leftarrow X$ . We denote  $C$  as the transformed representation of  $X$  in the low-dimensional manifold, whose base is inherited from  $D$ . This assumption can be rationalized as follows. 1)  $D \rightarrow C$ : a set of data points are usually embedded in a low-dimensional manifold. This finding can be dated back to the long history of dimensionality reduction [63, 53]. Nowadays, there are theoretical [3, 10] and empirical [82, 76] evidences showing that disentangled semantic manifolds emerge during training deep networks. 2)  $X \rightarrow C$ : features can be represented using (or projected onto) the manifold base linearly [64, 11] or non-linearly [8]. In particular, as later discussed in Section 3, we explicitly implement the base as feature dimensions (Figure 3(b)) and class-specific mean features (Figure 3(c)).

$X \rightarrow Y \leftarrow C$ . We denote  $Y$  as the classification effect (e.g., logits), which is determined by  $X$  via two ways: 1) the direct  $X \rightarrow Y$  and 2) the mediation  $X \rightarrow C \rightarrow Y$ . In particular, the first way can be removed if  $X$  can be fully represented by  $C$  (e.g., feature-wise adjustment in Section 3). The second way is inevitable even if the classifier does not take  $C$  as an explicit input, because any  $X$  can be

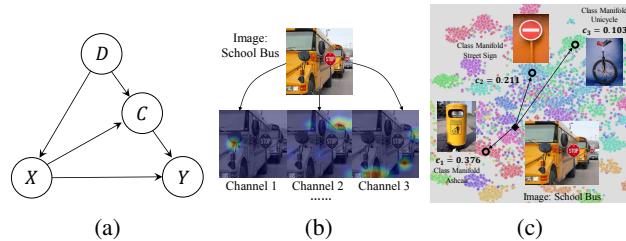


Figure 3: (a) Causal Graph for FSL; (b) Feature-wise illustration of  $D \rightarrow C$ : Feature channels of pre-trained network (e.g. 1 ... 512 for ResNet-10).  $X \rightarrow C$ : Per-channel response to an image (“school bus”) visualized by CAM[82]; (c) Class-wise illustration for  $D \rightarrow C$ : features are clustered according to the pre-training semantic classes (colored t-SNE plot[39]).  $X \rightarrow C$ : An image (“school bus”) can be represented in terms of the similarities among the base classes (“ashcan”, “unicycle”, “sign”).

I think this line  
is written as the  
“X” can contain  
informations which  
were not transferred  
to “C”, but effects “Y”

inherently represented by  $C$ . To illustrate, suppose that  $X$  is a linear combination of two base vectors plus a noise residual:  $\mathbf{x} = c_1 \mathbf{b}_1 + c_2 \mathbf{b}_2 + \mathbf{e}$ , any classifier  $f(\mathbf{x}) = f(c_1 \mathbf{b}_1 + c_2 \mathbf{b}_2 + \mathbf{e})$  will implicitly exploit the  $C$  representation in terms of  $\mathbf{b}_1$  and  $\mathbf{b}_2$ . In fact, this assumption also fundamentally validates unsupervised representation learning [7]. To see this, if  $C \not\rightarrow Y$  in Figure 3(a), uncovering the latent knowledge representation from  $P(Y|X)$  would be impossible, because the only path left that transfers knowledge from  $D$  to  $Y$ :  $D \rightarrow X \rightarrow Y$ , is cut off by conditioning on  $X$ :  $D \not\rightarrow X \rightarrow Y$ .

An ideal FSL model should capture the true causality between  $X$  and  $Y$  to generalize to unseen samples. For example, as illustrated in Figure 2(b), we expect that the “Lion” prediction is caused by the “lion” feature *per se*, but not the background “grass”. However, from the SCM in Figure 3(a), the conventional correlation  $P(Y|X)$  fails to do so, because the increased likelihood of  $Y$  given  $X$  is not only due to “ $X$  causes  $Y$ ” via  $X \rightarrow Y$  and  $X \rightarrow C \rightarrow Y$ , but also the spurious correlation via 1)  $D \rightarrow X$ , e.g., the “grass” knowledge generates the “grass” feature, and 2)  $D \rightarrow C \rightarrow Y$ , e.g., the “grass” knowledge generates the “grass” semantic, which provides useful context for “Lion” label. Therefore, to pursue the true causality between  $X$  and  $Y$ , we need to use the *causal intervention*  $P(Y|do(X))$  [48] instead of the likelihood  $P(Y|X)$  for the FSL objective.

### 2.3 Causal Intervention via Backdoor Adjustment

By now, an astute reader may notice that the causal graph in Figure 3(a) is also valid for Many-Shot Learning (MSL), *i.e.*, conventional learning based on pre-training. Compared to FSL, the  $P(Y|X)$  estimation of MSL is much more robust. For example, on *miniImageNet*, a 5-way-550-shot fine-tuned classifier can achieve 95% accuracy, while a 5-way-5-shot one only obtains 79%. We used to blame FSL for insufficient data by the law of large numbers in point estimation [16]. However, it does not answer why MSL converges to the true causal effects as the number of samples increases infinitely. In other words, why  $P(Y|do(X)) \approx P(Y|X)$  in MSL while  $P(Y|do(X)) \not\approx P(Y|X)$  in FSL?

To answer the question, we need to incorporate the endogenous feature sampling  $\mathbf{x} \sim P(X|I)$  into the estimation of  $P(Y|X)$ , where  $I$  denotes the sample ID. We have  $P(Y|X = \mathbf{x}_i) := \mathbb{E}_{\mathbf{x} \sim P(X|I)} P(Y|X = \mathbf{x}, I = i) = P(Y|I)$ , *i.e.*, we can use  $P(Y|I)$  to estimate  $P(Y|X)$ . In Figure 4(a), the causal relation between  $I$  and  $X$  is purely  $I \rightarrow X$ , *i.e.*,  $X \rightarrow I$  does not exist, because tracing the  $X$ ’s ID out of many-shot samples is like to find a needle in a haystack, given the nature that a DNN feature is an abstract and diversity-reduced representation of many samples [23]. However, as shown in Figure 4(b),  $X \rightarrow I$  persists in FSL, because it is much easier for a model to “guess” the correspondence, *e.g.*, the 1-shot extreme case that has a trivial 1-to-1 mapping for  $X \leftrightarrow I$ . Therefore, as we formally show in Appendix 1, the key causal difference between MSL and FSL is: **MSL essentially makes  $I$  an instrumental variable [1]** that achieves  $P(Y|X) := P(Y|I) \approx P(Y|do(X))$ . Intuitively, we can see that all the causalities between  $I$  and  $D$  in MSL are all blocked by colliders<sup>1</sup>, making  $I$  and  $D$  independent. So, the feature  $X$  is essentially “intervened” by  $I$ , no longer dictated by  $D$ , *e.g.*, neither “yellow grass” nor “green grass” dominates “Lion” in Figure 2(b), mimicking the causal intervention by controlling the use of pre-trained knowledge.

In this paper, we propose to use the backdoor adjustment [46] to achieve  $P(Y|do(X))$  without the need for many-shot, which certainly undermines the definition of FSL. The backdoor adjustment assumes that we can observe and stratify the confounder, *i.e.*,  $D = \{d\}$ , where each  $d$  is a stratification of the pre-trained knowledge. Formally, as shown in Appendix 2, the backdoor adjustment for the graph in Figure 3(a) is:

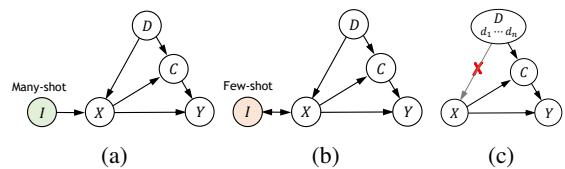


Figure 4: Causal graphs with sampling process. (a) Many-Shot Learning, where  $P(Y|X) \approx P(Y|do(X))$ ; (b) Few-Shot Learning where  $P(Y|X) \not\approx P(Y|do(X))$ ; (c) Interventional Few-Shot Learning where we directly model  $P(Y|do(X))$ .

$$P(Y|do(X = \mathbf{x})) = \sum_d P(Y|X = \mathbf{x}, D = d, C = g(\mathbf{x}, d)) P(D = d). \quad (1)$$

<sup>1</sup>In causal graph, the junction  $A \rightarrow B \leftarrow C$  is called a “collider”, making  $A$  and  $C$  independent even though  $A$  and  $C$  are linked via  $B$  [46]. For example,  $A$  = “Quality”,  $C$  = “Luck”, and  $B$  = “Paper Acceptance”.

where  $g$  is a function defined later. However, it is not trivial to instantiate  $d$ , especially when  $D$  is a 3rd-party delivered pre-trained network where the dataset is unobserved [22]. Next, we will offer three practical implementations of Eq. (1) for Interventional FSL.

### 3 Interventional Few-Shot Learning

Our implementation idea is inspired from the two inherent properties of any pre-trained DNN. First, each feature dimension carries a semantic meaning, e.g., every channel in convolutional neural network is well-known to encode visual concepts [82, 76]. So, each feature dimension represents a piece of knowledge. Second, most prevailing pre-trained models use a classification task as the objective, such as the 1,000-way classifier of ResNet [25] and the token predictor of BERT [17]. Therefore, the classifier can be considered as the distilled knowledge, which has been already widely adopted in literature [26]. Next, we will detail the proposed Interventional FSL (IFSL) by providing three different implementations<sup>2</sup> for  $g(\mathbf{x}, d)$ ,  $P(Y|X, D, C)$ , and  $P(D)$  in Eq. (1). In particular, the exact forms of  $P(Y|\cdot)$  across different classifiers are given in Appendix 5.

**Feature-wise Adjustment.** Suppose that  $\mathcal{F}$  is the index set of the feature dimensions of  $\mathbf{x}$ , e.g., from the last-layer of the pre-trained network  $\Omega$ . We divide  $\mathcal{F}$  into  $n$  equal-size disjoint subsets, e.g., the output feature dimension of ResNet-10 is 512, if  $n = 8$ , the  $i$ -th set will be a feature dimension index set of size  $512/8 = 64$ , i.e.,  $\mathcal{F}_i = \{64(i-1) + 1, \dots, 64i\}$ . The stratum set of pre-trained knowledge is defined as  $D := \{d_1, \dots, d_n\}$ , where each  $d_i = \mathcal{F}_i$ .

(i)  $g(\mathbf{x}, d_i) := \{k | k \in \mathcal{F}_i \cap \mathcal{I}_t\}$ , where  $\mathcal{I}_t$  is an index set whose corresponding absolute values in  $\mathbf{x}$  are larger than the threshold  $t$ . The reason is simple: if a feature dimension is inactive in  $\mathbf{x}$ , its corresponding adjustment can be omitted. We set  $t=1e-3$  in this paper.

(ii)  $P(Y|X, D, C) = P(Y|[\mathbf{x}]_c)$ , where  $c = g(\mathbf{x}, d_i)$  is implemented as the index set defined above,  $[\mathbf{x}]_c = \{\mathbf{x}_k\}_{k \in c}$  is a feature selector which selects the dimensions of  $\mathbf{x}$  according to the index set  $c$ . The classifier takes the adjusted feature  $[\mathbf{x}]_c$  as input. Note that  $d$  is already absorbed in  $c$ , so  $[\mathbf{x}]_c$  is essentially a function of  $(X, D, C)$ .

(iii)  $P(d_i) = 1/n$ , where we assume a uniform prior for the adjusted features.

(iv) The overall feature-wise adjustment is:

$$P(Y|do(X = \mathbf{x})) = \frac{1}{n} \sum_{i=1}^n P(Y|[\mathbf{x}]_c), \quad \text{where } c = \{k | k \in \mathcal{F}_i \cap \mathcal{I}_t\}. \quad (2)$$

It is worth noting that the feature-wise adjustment is always applicable, as we can always have the feature representation  $\mathbf{x}$  from the pre-trained network. Interestingly, our feature-wise adjustment sheds some light on the theoretical justifications for the multi-head trick in transformers [65]. We will explore this in future work.

**Class-wise Adjustment.** Suppose that there are  $m$  pre-training classes, denoted as  $\mathcal{A} = \{a_1, \dots, a_m\}$ . In class-wise adjustment, each stratum of pre-trained knowledge is defined as a pre-training class, i.e.,  $D := \{d_1, \dots, d_m\}$  and each  $d_i = a_i$ .

(i)  $g(\mathbf{x}, d_i) := P(a_i|\mathbf{x})\bar{\mathbf{x}}_i$ , where  $P(a_i|\mathbf{x})$  is the pre-trained classifier's probability output that  $\mathbf{x}$  belongs to class  $a_i$ , and  $\bar{\mathbf{x}}_i$  is the mean feature of pre-training samples from class  $a_i$ . Note that unlike feature-wise adjustment where  $c$  is an index set, here  $c = g(\mathbf{x}, d_i)$  is implemented as a real vector.

(ii)  $P(Y|X, D, C) = P(Y|\mathbf{x} \oplus g(\mathbf{x}, d_i))$ , where  $\oplus$  denotes vector concatenation.

(iii)  $P(d_i) = 1/m$ , where we assume a uniform prior of each class.

(iv) The overall class-wise adjustment is:

$$P(Y|do(X = \mathbf{x})) = \frac{1}{m} \sum_{i=1}^m P(Y|\mathbf{x} \oplus P(a_i|\mathbf{x})\bar{\mathbf{x}}_i) \approx P(Y|\mathbf{x} \oplus \frac{1}{m} \sum_{i=1}^m P(a_i|\mathbf{x})\bar{\mathbf{x}}_i), \quad (3)$$

---

<sup>2</sup>We assume that the combinations of the feature dimensions or classes are linear, otherwise the adjustment requires prohibitive  $\mathcal{O}(2^n)$  sampling. We will relax this assumption in future work.

where we adopt the Normalized Weighted Geometric Mean (NWGM) [71, 72] approximation to move the outer sum  $\sum P$  into the inner  $P(\sum)$ . This greatly reduces the network forward-pass consumption as  $m$  is usually large in pre-training dataset. Please refer to Appendix 3 for the detailed derivation.

**Combined Adjustment.** We can combine feature-wise and class-wise adjustment to make the stratification in backdoor adjustment much more fine-grained. Our combination is simple: applying feature-wise adjustment after class-wise adjustment. Thus, we have:

$$P(Y|do(X = \mathbf{x})) \approx \frac{1}{n} \sum_{i=1}^n P(Y|\mathbf{x})_c \oplus \frac{1}{m} \sum_{j=1}^m [P(a_j|\mathbf{x})\bar{\mathbf{x}}_j]_c, \text{ where } c = \{k | k \in \mathcal{F}_i \cap \mathcal{I}_t\}. \quad (4)$$

## 4 Related Work

**Few-Shot Learning.** FSL has a wide spectrum of methods, including fine-tuning [13, 18], optimizing model initialization [20, 42], generating model parameters [54, 36], learning a feature space for a better separation of sample categories [66, 77], feature transfer [58, 43], and transductive learning that additionally uses query set data [18, 29, 27]. Thanks to them, the classification accuracy has been drastically increased [29, 77, 73, 37]. However, accuracy as a single number cannot explain the paradoxical phenomenon in Figure 2. Our work offers an answer from a causal standpoint by showing that pre-training is a confounder. We not only further improve the accuracy of various FSL methods, but also explain the reason behind the improvements. In fact, the perspective offered by our work can benefit all the tasks that involve pre-training—any downstream task can be seen as FSL compared to the large-scale pre-training data.

**Negative Transfer.** The above phenomenon is also known as the negative transfer, where learning in source domain contributes negatively to the performance in target domain [44]. Many research works have been focused on when and how to conduct this transfer learning [30, 4, 81]. Yosinski *et al.* [74] split ImageNet according to man-made objects and natural objects as a test bed for feature transferability. They resemble the  $\mathcal{S} \not\propto \mathcal{Q}$  settings used in Figure 2(a). Other work also revealed that using deeper backbone might lead to degraded performance when the domain gap between training and test is large [33]. Some similar findings are reported in the few-shot setting [50] and NLP tasks [62]. Unfortunately, they didn't provide a theoretical explanation why it happens.

**Causal Inference.** Our work aims to deal with the pre-training confounder in FSL based on causal inference [48]. Causal inference was recently introduced to machine learning [40, 9] and has been applied to various fields in computer vision. [72] proposes a retrospective for image captioning and other applications include image classification [12, 38], imitation learning [15], long-tailed recognition [60] and semantic segmentation [78]. We are the first to approach FSL from a causal perspective. We would like to highlight that data-augmentation based FSL can also be considered as approximated intervention. These methods learn to generate additional support samples with image deformation [14, 79] or generative models [2, 80]. This can be viewed as physical interventions on the image features. Regarding the causal relation between image  $X$  and label  $Y$ , some works adopted anti-causal learning [41], *i.e.*,  $Y \rightarrow X$ , where the assumption is that labels  $Y$  are disentangled enough to be treated as Independent Mechanism (IM) [45, 59], which generates observed images  $X$  through  $Y \rightarrow X$ . However, our work targets at the more general case where labels can be entangled (*e.g.* “lion” and “dog” share the semantic “soft fur”) and the IM assumption may not hold. Therefore, we use causal prediction  $X \rightarrow Y$  as it is essentially a reasoning process, where the IM is captured by  $D$ , which is engineered to be disentangled through CNN (*e.g.*, the conv-operations are applied independently). In this way,  $D$  generates visual features through  $D \rightarrow X$  and emulates human's naming process through  $D \rightarrow Y$  (*e.g.*, “fur”, “four-legged”  $\rightarrow$  “meerkat”). In fact, the causal direction  $X \rightarrow Y$  (NOT anti-causal  $Y \rightarrow X$ ) has been empirically justified in complex CV tasks [32, 67, 60, 61].

## 5 Experiments

### 5.1 Datasets and Settings

**Datasets.** We conducted experiments on benchmark datasets in FSL literature: 1) **miniImageNet** [66] containing 600 images per class over 100 classes. We followed the split proposed in [51]: 64/16/20

classes for train/val/test. 2) ***tieredImageNet*** [52] is much larger compared to *miniImageNet* with 608 classes and each class around 1,300 samples. These classes were grouped into 34 higher-level concepts and then partitioned into 20/6/8 disjoint sets for train/val/test to achieve larger domain difference between training and testing. 3) Caltech-UCSD Birds-200-2011 (**CUB**) [70] for cross-domain evaluation. It contains 200 classes and each class has around 60 samples. The models used for CUB test were trained on the *miniImageNet*. Training and evaluation settings on *miniImageNet* and *tieredImageNet* are included in Appendix 5.

**Implementation Details.** We pre-trained the 10-layer ResNet (ResNet-10) [25] and the WideResNet (WRN-28-10) [75] as our backbones. Our proposed IFSL supports both fine-tuning and meta-learning. For fine-tuning, we applied average pooling on the last residual block and used the pooled features to train classifiers. For meta-learning, we deployed 5 representative methods that cover a large spectrum of meta-learning based FSL: 1) model initialization: MAML [20], 2) weight generator: LEO [54], transductive learning: SIB [29], 4) metric learning: MatchingNet (MN) [66], and 5) feature transfer: MTL [58]. For both fine-tuning and meta-learning, our IFSL aims to learn classifier  $P(Y|do(X))$  instead of the conventional  $P(Y|X)$ . Detailed implementations are given in Appendix 5.

**Evaluation Metrics.** Our evaluation is based on the following metrics: 1) Conventional accuracy (**Acc**) is the average classification accuracy commonly used in FSL [20, 66, 58]. 2) **Hardness-specific Acc**. For each query, we define a hardness that measures its semantic dissimilarity to the support set, and accuracy is then computed at different levels of query hardness. Specifically, query hardness is computed by  $h = \log((1 - s)/s)$  and  $s = \exp(\langle \mathbf{r}^+, \mathbf{p}_{c=gt}^+ \rangle / \sum_c \exp(\langle \mathbf{r}^+, \mathbf{p}_c^+ \rangle))$ , where  $\langle \cdot \rangle$  is the cosine similarity,  $(\cdot)^+$  represents the ReLU activation function,  $\mathbf{r}$  denotes the  $\Omega$  prediction logits of query,  $\mathbf{p}_c$  denotes the average prediction logits of class  $c$  in the support set and  $gt$  is the ground-truth of query. Using **Hardness-specific Acc** is similar to evaluating the hardness of FSL tasks [18], while ours is query-sample-specific and hence is more fine-grained. Later, we will show its effectiveness to unveil the spurious effects in FSL. 3) Feature localization accuracy (**CAM-Acc**) quantifies if a model “pays attention” to the actual object when making prediction. It is defined as the percentage of pixels inside the object bounding box by using Grad-CAM [56] score larger than 0.9. Compared to **Acc** that shows if the prediction is correct, **CAM-Acc** reveals whether the prediction is based on the correct visual cues.

Table 1: Acc (%) averaged over 2000 5-way FSL tasks before and after applying IFSL. We obtained the results by using official code and our backbones for a fair comparison across methods. We also implemented SIB in both transductive and inductive setting to facilitate fair comparison. For IFSL, we reported results of combined adjustment as it almost always outperformed feature-wise and class-wise adjustment. See Appendix 6 for Acc and 95% confidence intervals on all 3 types of adjustment.

Method	ResNet-10				WRN-28-10			
	<i>miniImageNet</i>		<i>tieredImageNet</i>		<i>miniImageNet</i>		<i>tieredImageNet</i>	
	5-shot	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot	1-shot
Linear	76.38	56.26	81.01	61.39	79.79	60.69	85.37	67.27
	+IFSL+ <b>2.19</b>	<b>77.97+1.59</b>	<b>60.13+3.87</b>	<b>82.08+1.07</b>	<b>64.29+2.9</b>	<b>80.97+1.18</b>	<b>64.12+3.43</b>	<b>86.19+0.82</b>
Cosine	76.68	56.40	81.13	62.08	79.72	60.83	85.41	67.30
	+IFSL+ <b>1.77</b>	<b>77.63+0.95</b>	<b>59.84+3.44</b>	<b>81.75+0.62</b>	<b>64.47+2.39</b>	<b>80.74+1.02</b>	<b>63.76+2.93</b>	<b>86.13+0.72</b>
<i>k</i> -NN	76.63	55.92	80.85	61.16	79.60	60.34	84.67	67.25
	+IFSL+ <b>3.13</b>	<b>78.42+1.79</b>	<b>62.31+6.36</b>	<b>81.98+1.13</b>	<b>65.71+4.55</b>	<b>81.08+1.48</b>	<b>64.98+4.64</b>	<b>86.06+1.39</b>
MAML [20]	70.85	56.59	74.02	59.17	73.92	58.02	77.20	61.40
	+IFSL+ <b>5.55</b>	<b>76.37+5.52</b>	<b>59.36+2.77</b>	<b>81.04+7.02</b>	<b>63.88+4.71</b>	<b>79.25+5.33</b>	<b>62.84+4.82</b>	<b>85.10+7.90</b>
LEO [54]	74.49	58.48	80.25	65.25	75.86	59.77	82.15	68.90
	+IFSL+ <b>1.94</b>	<b>76.91+2.42</b>	<b>61.09+2.61</b>	<b>81.43+1.18</b>	<b>66.03+0.78</b>	<b>77.72+1.86</b>	<b>62.19+2.42</b>	<b>85.04+2.89</b>
MTL [58]	75.65	58.49	81.14	64.29	77.30	62.99	83.23	70.08
	+IFSL+ <b>2.02</b>	<b>78.03+2.38</b>	<b>61.17+2.68</b>	<b>82.35+1.21</b>	<b>65.72+1.43</b>	<b>80.20+2.9</b>	<b>64.40+1.41</b>	<b>86.02+2.79</b>
MN [66]	75.21	61.05	79.92	66.01	77.15	63.45	82.43	70.38
	+IFSL+ <b>1.34</b>	<b>76.73+1.52</b>	<b>62.64+1.59</b>	<b>80.79+0.87</b>	<b>67.30+1.29</b>	<b>78.55+1.40</b>	<b>64.89+1.44</b>	<b>84.03+1.60</b>
SIB [29]	78.88	67.10	85.09	77.64	81.73	71.31	88.19	81.97
	(transductive)	+IFSL+ <b>1.15</b>	<b>80.32+1.44</b>	<b>68.85+1.75</b>	<b>85.43+0.34</b>	<b>78.03+0.39</b>	<b>83.21+1.48</b>	<b>73.51+2.20</b>
	(inductive)	+IFSL+ <b>2.05</b>	<b>77.68+2.04</b>	<b>60.33+3.13</b>	<b>82.75+1.06</b>	<b>67.34+1.83</b>	<b>80.05+1.88</b>	<b>63.14+3.02</b>

## 5.2 Results and Analysis

**Conventional Acc.** 1) From Table 1 we observe that IFSL consistently improves fine-tuning and meta-learning in all settings, which suggests that IFSL is agnostic to methods, datasets, and backbones. 2) In particular, the improvements are typically larger on 1-shot than 5-shot. For example, in fine-

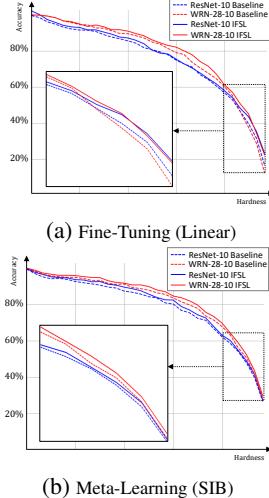


Figure 5: Accuracy across query hardness on 5-shot fine-tuning and meta-learning. Additional results are shown in Appendix 6.

Table 2: Comparison with state-of-the-arts of 5-way 1-/5-shot Acc (%) on *miniImageNet* and *tieredImageNet*.

Method	Backbone	<i>miniImageNet</i>		<i>tieredImageNet</i>	
		5-shot	1-shot	5-shot	1-shot
Baseline++ [12]	ResNet-10	75.90	53.97	-	-
IdeMe-Net <sup>†</sup> [14]	ResNet-10	73.78	57.61	80.34	60.32
TRAML [34]	ResNet-12	79.54	67.10	-	-
DeepEMD [77]	ResNet-12	82.41	65.91	86.03	71.16
CTM [35]	ResNet-18	80.51	64.12	84.28	68.41
FEAT [73]	WRN-28-10	81.80	66.69	84.38	70.41
Tran. Baseline [18]	WRN-28-10	78.40	65.73	85.50	73.34
wDAE-GNN [21]	WRN-28-10	78.85	62.96	83.09	68.18
SIB <sup>†</sup> [29]	WRN-28-10	81.73	71.31	88.19	81.97
<b>SIB+IFSL (ours)</b>	<b>WRN-28-10</b>	<b>83.21</b>	<b>73.51</b>	<b>88.69</b>	<b>83.07</b>

<sup>†</sup>Using our pre-trained backbone.

Table 3: Results of cross-domain evaluation: *miniImageNet* → CUB. The whole report is in Appendix 6.

Backbone	Method	5-shot	1-shot
		Linear	+IFSL
ResNet-10	Linear	58.84	42.25
	+IFSL	60.65	45.14
	SIB	60.60	45.87
	+IFSL	62.07	47.07
WRN-28-10	Linear	62.12	42.89
	+IFSL	64.15	45.64
	SIB	62.59	49.16
	+IFSL	64.43	50.71

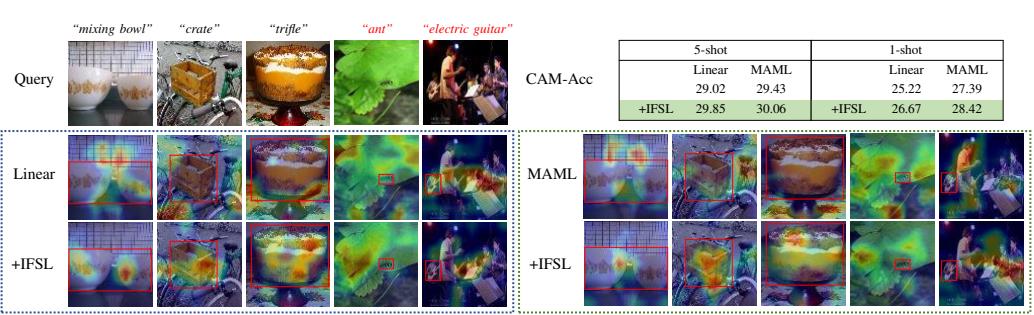


Figure 6: Some *miniImageNet* visualizations of Grad-Cam [56] activation of query images and the CAM-Acc (%) table of using linear classifier and MAML. Categories with red text represent failed cases. The complete results on CAM-Acc are shown in Appendix 6, where IFSL achieves similar or better results in all settings.

tuning, the average performance gain is 1.15% on 5-shot and 3.58% on 1-shot. The results support our analysis in Section 2.3 that FSL models are more prone to bias in lower-shot settings. 3) Regarding the average improvements on fine-tuning vs. meta-learning (*e.g.* *k*-NN and MN), we observe that IFSL improves more on fine-tuning in most cases. We conjecture that this is because meta-learning is an implicit form of intervention, where randomly sampled meta-training episodes effectively stratify the pre-trained knowledge. This suggests that meta-learning is fundamentally superior over fine-tuning due to increased robustness against confounders. We will investigate this potential theory in future work. 4) Additionally we see that the improvements on *miniImageNet* are usually larger than that on *tieredImageNet*. A possible reason is the much larger training set for *tieredImageNet*: it substantially increases the breadth of the pre-trained knowledge and the resulting models explain query samples much better. 5) According to Table 1 and Table 2, it is clear that our *k*-NN+IFSL outperforms IdeMe-Net [14] using the same pre-trained ResNet-10. This shows that using data augmentation — a method of physical data intervention as in IdeMe-Net [14] is inferior to our causal intervention in IFSL. 6) Overall, our IFSL achieves the new state-of-the-art on both datasets. Note that IFSL is flexible to be plugged into different baselines.

**Hardness-specific Acc.** 1) Figure 5(a) shows the plot of Hardness-specific Acc of fine-tuning. We notice that when query becomes harder, ResNet-10 (blue curves) becomes superior to WRN-28-10 (red curves). This tendency is consistent with Figure 2(a) illustrating the effect of the confounding bias caused by pre-training. 2) Intriguingly, in Figure 5(b) we notice that this tendency is reversed for meta-learning, *i.e.*, deeper backbone always performs better. The improved performance of deeper backbone on hard queries suggests that meta-learning should have some functions to remove the confounding bias. This evidence will inspire us to provide a causal view of meta-learning in future work. 3) Overall, Figure 5 shows that using IFSL further improves fine-tuning and meta-learning consistently across all hardness, validating the effectiveness of the proposed causal intervention.

**CAM-Acc & Visualization.** In Figure 6, we compare +IFSL to baseline linear classifier on the left and to baseline MAML [20] on the right, and summarize CAM-Acc results in the upper-right table. From the visualization, we see that using IFSL let the model pay more attention to the objects. However, notice that all models failed in the categories colored as red. A possible reason behind the failures is the extremely small size of the object — models have to resort to context for prediction. From the numbers, we can see our improvements for 1-shot are larger than that for 5-shot, consistent

with our findings using other evaluation metrics. These results suggest that IFSL helps models use the correct visual semantics for prediction by removing the confounding bias.

**Cross-Domain Generalization Ability.** In Table [3], we show the testing results on CUB using the models trained on the *miniImageNet*. The setting is challenging due to the big domain gap between the two datasets. We chose linear classifier as it outperforms cosine and  $k$ -NN in cross-domain setting and compared with transductive method — SIB. The results clearly show that IFSL works well in this setting and brings consistent improvements, with the average 1.94% of Acc. In addition, we can see that applying IFSL brings larger improvements to the inductive linear classifier than to the transductive SIB. It is possibly because transductive methods involve unlabeled query data and performs better than inductive methods with the additional information. Nonetheless we observe that IFSL can further improve SIB in cross-domain (Table [3]) and single-domain (Table [1]) generalization.

## 6 Conclusions

We presented a novel causal framework: Interventional Few-Shot Learning (IFSL), to address an overlooked deficiency in recent FSL methods: the pre-training is a confounder hurting the performance. Specifically, we proposed a structural causal model of the causalities in the process of FSL and then developed three practical implementations based on the backdoor adjustment. To better illustrate the deficiency, we diagnosed the classification accuracy comprehensively across query hardness, and showed that IFSL improves all the baselines across all the hardness. It is worth highlighting that the contribution of IFSL is not only about improving the performance of FSL, but also offering a causal explanation why IFSL works well: it is a causal approximation to many-shot learning. We believe that IFSL may shed light on exploring the new boundary of FSL, even though FSL is well-known to be ill-posed due to insufficient data. To upgrade IFSL, we will seek other observational intervention algorithms for better performance, and devise counterfactual reasoning for more general few-shot settings such as domain transfer.

## 7 Acknowledgements

The authors would like to thank all the anonymous reviewers for their constructive comments and suggestions. This research is partly supported by the Alibaba-NTU Singapore Joint Research Institute, Nanyang Technological University (NTU), Singapore; the Singapore Ministry of Education (MOE) Academic Research Fund (AcRF) Tier 1 and Tier 2 grant; and Alibaba Innovative Research (AIR) programme. We also want to thank Alibaba City Brain Group for the donations of GPUs.

## 8 Broader Impact

The proposed method aims to improve the Few-Shot Learning task. Advancements in FSL helps the deployment of machine learning models in areas where labelled data is difficult or expensive to obtain and it is closely related to social well-beings: few-shot drug discovery or medical imaging analysis in medical applications, cold-start item recommendation in e-commerce, few-shot reinforcement learning for industrial robots, etc.. Our method is based on causal inference and the analysis is rooted on causation rather than correlation. The marriage between causality and machine learning can produce more robust, transparent and explainable models, broadening the applicability of ML models and promoting fairness in artificial intelligence.

## References

- [1] Joshua D Angrist and Alan B Krueger. Instrumental variables and the search for identification: From supply and demand to natural experiments. *Journal of Economic Perspectives*, 2001. [4]
- [2] Antreas Antoniou, Amos Storkey, and Harrison Edwards. Data augmentation generative adversarial networks. In *Proceedings of the International Conference on Learning Representations Workshops*, 2018. [6]
- [3] Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix factorization. In *Advances in Neural Information Processing Systems*, 2019. [3]

- [4] Hossein Azizpour, Ali Sharif Razavian, Josephine Sullivan, Atsuto Maki, and Stefan Carlsson. Factors of transferability for a generic convnet representation. *IEEE transactions on Pattern Analysis and Machine Intelligence*, 2015. [6](#)
- [5] Pierre Baldi and Peter Sadowski. The dropout learning algorithm. *Artificial intelligence*, 2014. [16](#)
- [6] Alexander Balke and Judea Pearl. Bounds on treatment effects from studies with imperfect compliance. *Journal of the American Statistical Association*, 1997. [15](#)
- [7] Yoshua Bengio. Deep learning of representations for unsupervised and transfer learning. In *Proceedings of ICML Workshop on Unsupervised and Transfer Learning*, 2012. [4](#)
- [8] Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2013. [3](#)
- [9] Yoshua Bengio, Tristan Deleu, Nasim Rahaman, Rosemary Ke, Sébastien Lachapelle, Olexa Bilaniuk, Anirudh Goyal, and Christopher Pal. A meta-transfer objective for learning to disentangle causal mechanisms. In *International Conference on Learning Representations*, 2019. [6](#)
- [10] Michel Besserve, Rémy Sun, and Bernhard Schölkopf. Counterfactuals uncover the modular structure of deep generative models. *arXiv preprint arXiv:1812.03253*, 2018. [3](#)
- [11] Emmanuel J Candès, Xiaodong Li, Yi Ma, and John Wright. Robust principal component analysis? *Journal of the ACM (JACM)*, 2011. [3](#)
- [12] Krzysztof Chalupka, Pietro Perona, and Frederick Eberhardt. Visual causal feature learning. In *Uncertainty in Artificial Intelligence*, 2015. [6](#)
- [13] Wei-Yu Chen, Yen-Cheng Liu, Zsolt Kira, Yu-Chiang Frank Wang, and Jia-Bin Huang. A closer look at few-shot classification. In *International Conference on Learning Representations*, 2019. [1](#) [6](#) [8](#) [18](#) [19](#)
- [14] Zitian Chen, Yanwei Fu, Yu-Xiong Wang, Lin Ma, Wei Liu, and Martial Hebert. Image deformation meta-networks for one-shot learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. [6](#) [8](#)
- [15] Pim de Haan, Dinesh Jayaraman, and Sergey Levine. Causal confusion in imitation learning. In *Advances in Neural Information Processing Systems*, 2019. [6](#)
- [16] F.M. Dekking, C. Kraaikamp, H.P. Lopuhaä, and L.E. Meester. *A Modern Introduction to Probability and Statistics: Understanding Why and How*. Springer Texts in Statistics. Springer, 2005. [4](#)
- [17] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2019. [1](#) [5](#)
- [18] Guneet S Dhillon, Pratik Chaudhari, Avinash Ravichandran, and Stefano Soatto. A baseline for few-shot image classification. In *International Conference on Learning Representations*, 2020. [1](#) [6](#) [7](#) [8](#)
- [19] Li Fei-Fei, Rob Fergus, and Pietro Perona. One-shot learning of object categories. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2006. [1](#)
- [20] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International Conference on Machine Learning*, 2017. [2](#) [3](#) [6](#) [7](#) [8](#) [17](#) [18](#) [20](#) [21](#) [22](#) [23](#)
- [21] Spyros Gidaris and Nikos Komodakis. Generating classification weights with gnn denoising autoencoders for few-shot learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. [2](#) [3](#) [8](#)
- [22] Ross Girshick, Ilya Radosavovic, Georgia Gkioxari, Piotr Dollár, and Kaiming He. Detectron. <https://github.com/facebookresearch/detectron>, 2018. [5](#)
- [23] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. The MIT Press, 2016. [4](#)
- [24] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE International Conference on Computer Vision*, 2017. [1](#)

- [25] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016. [1, 5, 7, 17, 18]
- [26] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. In *Advances in Neural Information Processing Systems Deep Learning Workshop*, 2014. [1, 5]
- [27] Ruibing Hou, Hong Chang, MA Bingpeng, Shiguang Shan, and Xilin Chen. Cross attention network for few-shot classification. In *Advances in Neural Information Processing Systems*, 2019. [6]
- [28] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017. [1]
- [29] Shell Xu Hu, Pablo Moreno, Yang Xiao, Xi Shen, Guillaume Obozinski, Neil Lawrence, and Andreas Damianou. Empirical bayes transductive meta-learning with synthetic gradients. In *International Conference on Learning Representations*, 2020. [2, 6, 7, 8, 19, 20, 21, 22, 23]
- [30] Minyoung Huh, Pulkit Agrawal, and Alexei A Efros. What makes imagenet good for transfer learning? *arXiv preprint arXiv:1608.08614*, 2016. [6]
- [31] Muhammad Abdullah Jamal and Guo-Jun Qi. Task agnostic meta-learning for few-shot learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. [1]
- [32] Qi Jiaxin, Niu Yulei, Huang Jianqiang, and Zhang Hanwang. Two causal principles for improving visual dialog. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020. [6]
- [33] Simon Kornblith, Jonathon Shlens, and Quoc V Le. Do better imagenet models transfer better? In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2019. [6]
- [34] Aoxue Li, Weiran Huang, Xu Lan, Jiashi Feng, Zhenguo Li, and Liwei Wang. Boosting few-shot learning with adaptive margin loss. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020. [8]
- [35] Hongyang Li, David Eigen, Samuel Dodge, Matthew Zeiler, and Xiaogang Wang. Finding task-relevant features for few-shot learning by category traversal. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. [8]
- [36] Huaiyu Li, Weiming Dong, Xing Mei, Chongyang Ma, Feiyue Huang, and Bao-Gang Hu. Lgm-net: Learning to generate matching networks for few-shot learning. In *International Conference on Machine Learning*, 2019. [6]
- [37] Yaoyao Liu, Bernt Schiele, and Qianru Sun. An ensemble of epoch-wise empirical bayes for few-shot learning. In *European Conference on Computer Vision*, 2020. [6]
- [38] David Lopez-Paz, Robert Nishihara, Soumith Chintala, Bernhard Scholkopf, and Léon Bottou. Discovering causal signals in images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017. [6]
- [39] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research*, 2008. [3]
- [40] Sara Magliacane, Thijs van Ommen, Tom Claassen, Stephan Bongers, Philip Versteeg, and Joris M Mooij. Domain adaptation by using causal inference to predict invariant conditional distributions. In *Advances in Neural Information Processing Systems*, 2018. [6]
- [41] Gong Mingming, Zhang Kun, Liu Tongliang, Tao Dacheng, Clark Glymour, and Bernhard Schölkopf. Domain adaptation with conditional transferable components. In *International Conference on Machine Learning*, 2016. [6]
- [42] Alex Nichol and John Schulman. Reptile: a scalable metalearning algorithm. *arXiv preprint arXiv:1803.02999*, 2018. [6]
- [43] Boris Oreshkin, Pau Rodríguez López, and Alexandre Lacoste. Tadam: Task dependent adaptive metric for improved few-shot learning. In *Advances in Neural Information Processing Systems*, 2018. [6]
- [44] Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on Knowledge and Data Engineering*, 2009. [6]

- [45] Giambattista Parascandolo, Niki Kilbertus, Mateo Rojas-Carulla, and Bernhard Schölkopf. Learning independent causal mechanisms. In *International Conference on Machine Learning*, 2018. [6](#)
- [46] J. Pearl, M. Glymour, and N.P. Jewell. *Causal Inference in Statistics: A Primer*. Wiley, 2016. [2](#)  
[3](#) [4](#) [14](#)
- [47] Judea Pearl. Causal diagrams for empirical research. *Biometrika*, 1995. [15](#)
- [48] Judea Pearl. *Causality: Models, Reasoning and Inference*. Cambridge University Press, 2nd edition, 2009. [4](#) [6](#) [14](#)
- [49] Siyuan Qiao, Chenxi Liu, Wei Shen, and Alan L Yuille. Few-shot image recognition by predicting parameters from activations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018. [2](#)
- [50] Tiago Ramalho, Thierry Soubie, and Stefano Peluchetti. An empirical study of pretrained representations for few-shot classification. *arXiv preprint arXiv:1910.01319*, 2019. [6](#)
- [51] Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. In *International Conference on Learning Representations*, 2017. [6](#)
- [52] Mengye Ren, Eleni Triantafillou, Sachin Ravi, Jake Snell, Kevin Swersky, Joshua B. Tenenbaum, Hugo Larochelle, and Richard S. Zemel. Meta-learning for semi-supervised few-shot classification. In *International Conference on Learning Representations*, 2018. [2](#) [7](#)
- [53] Sam T Roweis and Lawrence K Saul. Nonlinear dimensionality reduction by locally linear embedding. *Science*, 2000. [3](#)
- [54] Andrei A. Rusu, Dushyant Rao, Jakub Sygnowski, Oriol Vinyals, Razvan Pascanu, Simon Osindero, and Raia Hadsell. Meta-learning with latent embedding optimization. In *International Conference on Learning Representations*, 2019. [6](#) [7](#) [17](#) [19](#) [20](#) [21](#) [22](#) [23](#)
- [55] Adam Santoro, Sergey Bartunov, Matthew Botvinick, Daan Wierstra, and Timothy Lillicrap. Meta-learning with memory-augmented neural networks. In *International Conference on Machine Learning*, 2016. [1](#)
- [56] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE International Conference on Computer Vision*, 2017. [7](#) [8](#)
- [57] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In *Advances in Neural Information Processing Systems*, 2017. [16](#) [18](#)
- [58] Qianru Sun, Yaoyao Liu, Tat-Seng Chua, and Bernt Schiele. Meta-transfer learning for few-shot learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. [6](#) [7](#) [18](#) [19](#) [20](#) [21](#) [22](#) [23](#)
- [59] Raphael Suter, Đorđe Miladinović, Bernhard Schölkopf, and Stefan Bauer. Robustly disentangled causal mechanisms: Validating deep representations for interventional robustness. In *International Conference on Machine Learning*, 2019. [6](#)
- [60] Kaihua Tang, Jianqiang Huang, and Hanwang Zhang. Long-tailed classification by keeping the good and removing the bad momentum causal effect. In *Advances in Neural Information Processing Systems*, 2020. [6](#)
- [61] Kaihua Tang, Yulei Niu, Jianqiang Huang, Jiaxin Shi, and Hanwang Zhang. Unbiased scene graph generation from biased training. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020. [6](#)
- [62] Marc Tanti, Albert Gatt, and Kenneth P Camilleri. Transfer learning from language models to image caption generators: Better models may not transfer better. *arXiv preprint arXiv:1901.01216*, 2019. [6](#)
- [63] Joshua B Tenenbaum, Vin De Silva, and John C Langford. A global geometric framework for nonlinear dimensionality reduction. *Science*, 2000. [3](#)
- [64] Matthew Turk and Alex Pentland. Face recognition using eigenfaces. In *Proceedings. 1991 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 1991. [3](#)

- [65] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, 2017. [1](#) [5](#)
- [66] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for one shot learning. In *Advances in Neural Information Processing Systems*, 2016. [2](#) [3](#) [6](#) [7](#) [19](#) [20](#) [21](#) [22](#) [23](#)
- [67] Tan Wang, Jianqiang Huang, Hanwang Zhang, and Qianru Sun. Visual commonsense r-cnn. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. [6](#)
- [68] Yan Wang, Wei-Lun Chao, Kilian Q Weinberger, and Laurens van der Maaten. Simpleshot: Revisiting nearest-neighbor classification for few-shot learning. *arXiv preprint arXiv:1911.04623*, 2019. [18](#)
- [69] Yaqing Wang and Quanming Yao. Few-shot learning: A survey. *arXiv preprint arXiv:1904.05046*, 2019. [1](#)
- [70] P. Welinder, S. Branson, T. Mita, C. Wah, F. Schroff, S. Belongie, and P. Perona. Caltech-UCSD Birds 200. Technical report, California Institute of Technology, 2010. [2](#) [7](#)
- [71] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *International Conference on Machine Learning*, 2015. [6](#)
- [72] Xu Yang, Hanwang Zhang, and Jianfei Cai. Deconfounded image captioning: A causal retrospect. *arXiv preprint arXiv:2003.03923*, 2020. [6](#)
- [73] Han-Jia Ye, Hexiang Hu, De-Chuan Zhan, and Fei Sha. Few-shot learning via embedding adaptation with set-to-set functions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020. [6](#) [8](#)
- [74] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural networks? In *Advances in Neural Information Processing Systems*, 2014. [6](#)
- [75] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. In *British Machine Vision Conference*, 2016. [7](#) [17](#) [18](#)
- [76] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *European Conference on Computer Vision*, 2014. [3](#) [5](#)
- [77] Chi Zhang, Yujun Cai, Guosheng Lin, and Chunhua Shen. Deepemd: Few-shot image classification with differentiable earth mover’s distance and structured classifiers. 2020. [6](#) [8](#)
- [78] Dong Zhang, Hanwang Zhang, Jinhui Tang, Xian sheng Hua, and Qianru Sun. Causal intervention for weakly-supervised semantic segmentation. In *Advances in Neural Information Processing Systems*, 2020. [6](#)
- [79] Hongguang Zhang, Jing Zhang, and Piotr Koniusz. Few-shot learning via saliency-guided hallucination of samples. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. [6](#)
- [80] Ruixiang Zhang, Tong Che, Zoubin Ghahramani, Yoshua Bengio, and Yangqiu Song. Metagan: An adversarial approach to few-shot learning. In *Advances in Neural Information Processing Systems*, 2018. [6](#)
- [81] Youshan Zhang and Brian D Davison. Impact of imagenet model selection on domain adaptation. In *Proceedings of the IEEE Winter Conference on Applications of Computer Vision Workshops*, 2020. [6](#)
- [82] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016. [3](#) [5](#)

---

# Supplementary Material for Interventional Few-Shot Learning

---

Zhongqi Yue<sup>1,3</sup>, Hanwang Zhang<sup>1</sup>, Qianru Sun<sup>2</sup>, Xian-Sheng Hua<sup>3</sup>

<sup>1</sup>Nanyang Technological University, <sup>2</sup>Singapore Management University,

<sup>3</sup>Damo Academy, Alibaba Group

yuez0003@ntu.edu.sg, hanwangzhang@ntu.edu.sg,  
qianrusun@smu.edu.sg, xiansheng.hxs@alibaba-inc.com

This supplementary material is organized as follows:

- Section A.1 details our analysis in Section 2.3 by showing many-shot learning converges to true causal effect through instrumental variable (IV);
- Section A.2 gives the derivation for the backdoor adjustment formula in Eq. (1);
- Section A.3 presents the detailed derivation for the NWGM approximation used in Eq. (3) and (4);
- Section A.4 includes the algorithms for adding IFSL to fine-tuning and meta-learning;
- Section A.5 shows the implementation details for pre-training (Section A.5.1), fine-tuning (Section A.5.2) and meta-learning (Section A.5.3);
- Section A.6 includes additional experimental results on Conventional Acc (Section A.6.1), Hardness-Specific Acc (Section A.6.2), CAM-Acc (Section A.6.3) and cross-domain evaluation (Section A.6.4).

## A.1 Instrumental Variable

In this section, we will show that in our causal graph for many-shot learning, the sampling ID  $I$  is essentially an instrumental variable for  $X \rightarrow Y$  that achieves  $P(Y|I) \approx P(Y|do(X))$ . Before introducing instrumental variable, we first formally define *d-separation* [46], which gives a criterion to study the dependencies between nodes (data variables) in any structural causal model.

**d-separation.** A set of nodes  $Z$  blocks a path  $p$  if and only if 1)  $p$  contains a *chain*  $A \rightarrow B \rightarrow C$  or a *fork*  $A \leftarrow B \rightarrow C$  and the middle node  $B$  is in  $Z$ ; 2)  $p$  contains a *collider*  $A \rightarrow B \leftarrow C$  such that the middle node  $B$  and its descendants are not in  $Z$ . If conditioning on  $Z$  blocks every path between  $X$  and  $Y$ , we say  $X$  and  $Y$  are *d-separated* conditional on  $Z$ , i.e.,  $X$  and  $Y$  are independent given  $Z$  ( $X \perp\!\!\!\perp Y|Z$ ).

**Instrumental Variable.** For a structural causal model  $\mathcal{G}$ , a variable  $Z$  is an *instrumental variable* (IV) to  $X \rightarrow Y$  by satisfying the graphical criteria [48]: 1)  $(Z \perp\!\!\!\perp Y)_{\mathcal{G}_{\overline{X}}}$ ; 2)  $(Z \not\perp\!\!\!\perp X)_{\mathcal{G}}$ , where  $\mathcal{G}_{\overline{X}}$  is the manipulated graph where all incoming arrows to node  $X$  are deleted. For the SCM of many-shot learning in Figure 4(a), it is easy to see that  $I$  satisfies both criteria and therefore it is an IV for  $X \rightarrow Y$ . However, in the few-shot SCM in Figure 4(b), the paths  $I \leftarrow X \leftarrow D \rightarrow C \rightarrow Y$  and  $I \leftarrow X \rightarrow C \rightarrow Y$  are not blocked in  $\mathcal{G}_{\overline{X}}$ , which means the first criterion is not met  $(I \not\perp\!\!\!\perp Y)_{\mathcal{G}_{\overline{X}}}$  and  $I$  is not an instrumental variable in the few-shot learning case.

Instrumental variable can help find the true causal effect even in the presence of confounder. This is due to the collider junction that makes the IV and confounder independent ( $I \perp\!\!\!\perp D$  in Figure 4(a)). To see this, we will first consider a simplified case of Figure 4(a) where each causal link represents a linear relationship and we aim to find the true causal effect from  $X \rightarrow Y$  through linear regression. Without loss of generality, let  $I, X, Y$  take the value of real number. Denote  $r_{IX}, r_{XY}$ ,

and  $r_{IY}$  as the slope of regression line between  $I$  and  $X$ ,  $X$  and  $Y$ ,  $I$  and  $Y$  respectively. Notice that  $r_{XY}$  is spurious as it is contaminated by the backdoor path  $X \leftarrow D \rightarrow C \rightarrow Y$ . However, since the path  $I \rightarrow X \leftarrow D \rightarrow C \rightarrow Y$  is blocked due to collider at  $X$ ,  $r_{IY}$  is free from confounding bias. Therefore  $r_{IY}/r_{IX}$  gives the true causal effect from  $X \rightarrow Y$ . Similarly, in the classification case of many-shot learning, a classifier is trained to maximize the conditional probability on the IV  $P(Y|I)$ . As the ID-sample matching  $I \rightarrow X$  is deterministic, the classifier eventually learns to predict based on the true causal relationship  $X \rightarrow Y$ . Yet in the complex case of image classification, it is unreasonable to assume linear relationships between variables. In the nonlinear case, it is shown in [6] that observations on IV provide a bound for the true causal effect. This means that learning based on  $P(Y|I)$  provides an approximation to the true causal effect, *i.e.*  $P(Y|I) \approx P(Y|do(X))$ .

## A.2 Derivation of Backdoor Adjustment for the Proposed Causal Graph

We will show the derivation of the backdoor adjustment for the causal graph in Figure 3(a) using the three rules of *do*-calculus [47].

For a causal directed acyclic graph  $\mathcal{G}$ , let  $X, Y, Z$  and  $W$  be arbitrary disjoint sets of nodes. We use  $\mathcal{G}_{\bar{X}}$  to denote the manipulated graph where all incoming arrows to node  $X$  are deleted. Similarly  $\mathcal{G}_{\underline{X}}$  represents the graph where outgoing arrows from node  $X$  are deleted. We use lower case  $x, y, z$  and  $w$  for specific values taken by each set of nodes:  $X = x, Y = y, Z = z$  and  $W = w$ . For any interventional distribution compatible with  $\mathcal{G}$ , we have the following three rules:

**Rule 1** Insertion/deletion of observations:

$$P(y|do(x), z, w) = P(y|do(x), w), \text{if}(Y \perp\!\!\!\perp Z|X, W)_{\mathcal{G}_{\bar{X}}} \quad (\text{A5})$$

**Rule 2** Action/observation exchange:

$$P(y|do(x), do(z), w) = P(y|do(x), z, w), \text{if}(Y \perp\!\!\!\perp Z|X, W)_{\mathcal{G}_{\bar{X}Z}} \quad (\text{A6})$$

**Rule 3** Insertion/deletion of actions:

$$P(y|do(x), do(z), w) = P(y|do(x), w), \text{if}(Y \perp\!\!\!\perp Z|X, W)_{\mathcal{G}_{\bar{X}Z(W)}}, \quad (\text{A7})$$

where  $Z(W)$  is the set of nodes in  $Z$  that are not ancestors of any  $W$ -node in  $\mathcal{G}_{\bar{X}}$ .

In our causal graph, the desired interventional distribution  $P(Y|do(X = \mathbf{x}))$  can be derived by:

$$P(Y|do(\mathbf{x})) = \sum_d P(Y|do(X = \mathbf{x}), D = d)P(D = d|do(X = \mathbf{x})) \quad (\text{A8})$$

$$= \sum_d P(Y|do(X = \mathbf{x}), D = d)P(D = d) \quad (\text{A9})$$

$$= \sum_d P(Y|X = \mathbf{x}, D = d)P(D = d) \quad (\text{A10})$$

$$= \sum_d \sum_c P(Y|X = \mathbf{x}, D = d, C = c)P(C = c|X = \mathbf{x}, D = d)P(D = d) \quad (\text{A11})$$

$$= \sum_d P(Y|X = \mathbf{x}, D = d, C = g(\mathbf{x}, d))P(D = d), \quad (\text{A12})$$

where Eq. (A8) and Eq. (A11) follow the law of total probability; Eq. (A9) uses Rule 3 given  $D \perp\!\!\!\perp X$  in  $\mathcal{G}_{\bar{X}}$ ; Eq. (A10) uses Rule 2 to change the intervention term to observation as  $(Y \perp\!\!\!\perp X|D)$  in  $\mathcal{G}_{\underline{X}}$ ; Eq. (A12) is because in our causal graph,  $C$  takes a deterministic value given by function  $g(\mathbf{x}, d)$ . This reduces summation over all values of  $C$  in Eq. (A11) to a single probability measure in Eq. (A12).

## A.3 Derivation of NWGM Approximation

We will show the derivation of NWGM approximation used in Eq. (3) and (4). In a  $K$ -way FSL problem, let  $f(\cdot)$  be a classifier function that calculates logits for  $K$  classes and  $\sigma$  be the softmax

function over  $K$  classes. The approximation effectively moves the outer expectation inside the classifier function:  $\mathbb{E}[\sigma(f(\cdot))] \approx \sigma(f(\mathbb{E}[\cdot]))$ .

We will first show the derivation for moving the expectation inside softmax function, *i.e.*,  $\mathbb{E}[\sigma(f(\cdot))] \approx \sigma(\mathbb{E}[f(\cdot)])$ . Without loss of generality, the backdoor adjustment formula in Eq. (3) and Eq. (4) can be written as:

$$P(Y = y | do(X = \mathbf{x})) = \sum_{d \in D} \sigma(f_y(\mathbf{x} \oplus \mathbf{c})) P(d), \quad (\text{A13})$$

where  $D$  represents the set of stratifications,  $f_y$  is the classifier logit for class  $y$ ,  $\mathbf{c} = g(\mathbf{x}, d)$  is the feature concatenated to  $\mathbf{x}$  in Eq. (3) and (4) and  $P(d)$  is the prior for each stratification.

It is shown in [5] that Eq. (A13) can be approximated by the Normalized Weighted Geometric Mean (NWGM) as:

$$\sum_{d \in D} \sigma(f_y(\mathbf{x} \oplus \mathbf{c})) P(d) \approx NWGM_{d \in D}(\sigma(f_y(\mathbf{x} \oplus \mathbf{c}))) \quad (\text{A14})$$

$$= \frac{\prod_d [\exp(f_y(\mathbf{x} \oplus \mathbf{c}))]^{P(d)}}{\sum_{i=1}^K \prod_d [\exp(f_i(\mathbf{x} \oplus \mathbf{c}))]^{P(d)}} \quad (\text{A15})$$

$$= \frac{\exp(\sum_d f_y(\mathbf{x} \oplus \mathbf{c}) P(d))}{\sum_{i=1}^K \exp(\sum_d f_i(\mathbf{x} \oplus \mathbf{c}) P(d))} \quad (\text{A16})$$

$$= \sigma(\mathbb{E}_d[f_y(\mathbf{x} \oplus \mathbf{c})]), \quad (\text{A17})$$

where Eq. (A14) follows [5], Eq. (A15) follows the definition of NWGM, Eq. (A16) is because  $\exp(a)^b = \exp(ab)$ .

Next we will show the derivation for linear, cosine and  $k$ -NN classifier to further move the expectation inside the classifier function, *i.e.*,  $\sigma(\mathbb{E}[f(\cdot)]) = \sigma(f(\mathbb{E}[\cdot]))$ .

For the linear classifier,  $f(\mathbf{x} \oplus \mathbf{c}) = \mathbf{W}_1 \mathbf{x} + \mathbf{W}_2 \mathbf{c}$ , where  $\mathbf{W}_1, \mathbf{W}_2 \in \mathbb{R}^{K \times N}$  denote the learnable weight,  $N$  is the feature dimension, which is the same for  $\mathbf{x}$  and  $\mathbf{c}$  in Eq. (3) and (4). The bias term is dropped as it does not impact our analysis. Now the expectation can be further moved inside the classifier function through:

$$\sum_d f(\mathbf{x} \oplus \mathbf{c}) P(d) = \sum_d (\mathbf{W}_1 \mathbf{x} + \mathbf{W}_2 \mathbf{c}) P(d) \quad (\text{A18})$$

$$= \mathbf{W}_1 \mathbf{x} + \sum_d \mathbf{W}_2 \mathbf{c} P(d) \quad (\text{A19})$$

$$= f(\mathbf{x} \oplus \sum_d \mathbf{c} P(d)), \quad (\text{A20})$$

where Eq. (A19) is because the feature vector  $\mathbf{x}$  is the same for all  $d$  and  $\mathbb{E}_d[\mathbf{x}] = \mathbf{x}$ .

For the cosine classifier,  $f(\mathbf{x} \oplus \mathbf{c}) = (\mathbf{W}_1 \mathbf{x} + \mathbf{W}_2 \mathbf{c}) / \|\mathbf{x} \oplus \mathbf{c}\| \|\mathbf{W}\|$ , where  $\mathbf{W} \in \mathbb{R}^{K \times 2N}$  is the concatenation of  $\mathbf{W}_1$  and  $\mathbf{W}_2$ . In the special case where  $\mathbf{x}$  and  $\mathbf{c}$  are unit vector,  $\|\mathbf{x} \oplus \mathbf{c}\|$  is  $\sqrt{2}$  and the cosine classifier function reduces to a linear combination of terms involving only  $\mathbf{x}$  and only  $\mathbf{c}$ . From there, the analysis for linear classifier follows and we have  $\sigma(\mathbb{E} f(\cdot)) = \sigma(f(\mathbb{E} \cdot))$  for cosine classifier. In the general case where  $\mathbf{x}$  and  $\mathbf{c}$  are not unit vector, moving the expectation inside cosine classifier function is an approximation  $\sigma(\mathbb{E}[f(\cdot)]) \approx \sigma(f(\mathbb{E}[\cdot]))$ .

For the  $k$ -NN classifier, our implementation calculates class centroids using the mean feature of the  $K$  support sets and then uses the nearest centroid for prediction (1-NN). Specifically, let  $\mathbf{x}$  be a feature vector and  $\mathbf{x}'$  be the  $i$ th class centroid,  $i \in \{1, \dots, K\}$ . The logit for class  $i$  is given by  $f_i(\mathbf{x}) = -\|\mathbf{x} - \mathbf{x}'\|^2$ . It is shown in [57] that  $k$ -NN classifier that uses squared Euclidean distance to generate logits is equivalent to a linear classifier with a particular parameterization. Therefore, our analysis on linear classifier follows for  $k$ -NN.

In summary, the derivation of  $\mathbb{E}[\sigma(f(\cdot))] \approx \sigma(f(\mathbb{E}[\cdot]))$  is a two-stage process where we first move the expectation inside the softmax function and then further move it inside the classifier function.

## A.4 Algorithms for Fine-tuning and Meta-Learning with IFSL

In this section, we will briefly revisit the settings of fine-tuning and meta-learning and introduce how to integrate IFSL into them.

In fine-tuning, the goal is to train a classifier  $\theta$  conditioned on the current support set  $\mathcal{S} = \{(\mathbf{x}_i, y_i)\}_{i=1}^{n_s}$ , where  $\mathbf{x}_i$  is the feature generated by  $\Omega$  for  $i$ th sample,  $y_i$  is the ground-truth label for  $i$ th sample and  $n_s$  is the support set size. This is achieved by first predicting the support label  $\hat{y}$  using the classifier  $P(y|\mathbf{x}; \theta)$ . Then with the predicted label  $\hat{y}$  and ground-truth label  $y$ , one can calculate a loss  $\mathcal{L}(\hat{y}, y)$  (usually cross-entropy loss) to update the classifier parameter, e.g. through stochastic gradient descent. Adding IFSL to fine-tuning is simple: 1) Pick an adjustment strategy introduced in Section 3. Each implementation defines the set of pre-trained knowledge stratifications  $D$ , function form of  $g(X, D)$ , function form of  $P(Y|X, D, C)$  and the prior  $P(D)$ ; 2) The classifier prediction is now based on  $P(Y|do(X); \theta)$ . The process of fine-tuning with IFSL is summarized in Algorithm 1. Note that for the non-parametric  $k$ -NN classifier, the fine-tuning process is not applicable. When adding IFSL to  $k$ -NN, each sample is represented by the *adjusted feature* instead of original feature  $\mathbf{x}$ . Please refer to the classifier inputs in Eq. (2), (3) and (4) for the exact form of adjusted feature.

In meta-learning, the goal is to learn the additional “learning behavior” parameterized by  $\phi$  using training episodes  $\{(\mathcal{S}_i, \mathcal{Q}_i)\}$  sampled from training dataset  $\mathcal{D}$ . The classifier in meta-learning makes predictions by additionally conditioning on the learning behavior, written as  $P_\phi(y|\mathbf{x}; \theta)$ . Within each episode,  $\theta$  is first fine-tuned on the support set  $\mathcal{S}_i$ . Then the fine-tuned model is tested on the query set  $\mathcal{Q}_i$  to obtain the loss  $\mathcal{L}_\phi(\mathcal{S}_i, \mathcal{Q}_i)$  (e.g. using cross-entropy loss). Finally the loss is used to update  $\phi$  using an optimizer. It is also easy to integrate IFSL into meta-learning by only changing the classifier from  $P_\phi(y|\mathbf{x}; \theta)$  to  $P_\phi(y|do(\mathbf{x}); \theta)$ . The flow of meta-learning with IFSL is presented in Algorithm 2. Firstly notice that the initialization of  $\theta$  in each task may depend on  $\phi$  or  $\mathcal{S}_i$ . For example, in MAML [20]  $\phi$  essentially defines an initialization of model parameters, and in LEO [54] the initial classifier parameter is generated conditioned on  $\phi$  and  $\mathcal{S}_i$ . Secondly, although the fine-tuning of  $\theta$  largely follows Algorithm 1, some meta-learning methods additionally utilize meta-knowledge  $\phi$ . For example, in SIB the gradients for updating  $\theta$  are predicted by  $\phi$  using unlabelled query features instead of calculated from  $\mathcal{L}(\hat{y}, y)$  as in Algorithm 1.

---

### Algorithm 1: Fine-tuning + IFSL

```

Input :  $D$ , Support set  $\mathcal{S} = \{(\mathbf{x}_i, y_i)\}_{i=1}^{n_s}$ 
Output : Fine-tuned classifier parameters  $\theta$ 
Initialize  $\theta$ ;
while not converged do
    for  $i = 1, \dots, n_s$  do
        for  $d \in D$  do
            Calculate  $c = g(\mathbf{x}, d)$ ;
            Obtain  $P(Y|\mathbf{x}_i, c, d; \theta), P(d)$ 
            Prediction  $\hat{y}_i = P(y|do(\mathbf{x}); \theta)$ ;
            Update  $\theta$  using  $\mathcal{L}(\hat{y}_i, y_i)$ 
    return  $\theta$ 

```

---



---

### Algorithm 2: Meta-Learning + IFSL

```

Input :  $D$ , training dataset  $\mathcal{D}$ 
Output : Optimized meta-parameters  $\phi$ 
Initialize  $\phi$ ;
while not converged do
    Sample  $(\mathcal{S}_i, \mathcal{Q}_i)$  from  $\mathcal{D}$  ;
    Initialize classifier  $\theta$  with  $\phi, \mathcal{S}_i$ ;
    Fine-tune  $\theta$  using Algorithm 1
        conditioned on  $\phi$  ;
    Predict query based on  $P_\phi(y|do(\mathbf{x}); \theta)$ ;
    Update  $\phi$  using  $\mathcal{L}_\phi(\mathcal{S}_i, \mathcal{Q}_i; \theta)$ 
return  $\phi$ 

```

---

## A.5 Implementation Details

### A.5.1 Pre-training

Prior to fine-tuning or meta-learning, we pre-trained a deep neural network (DNN) as feature extractor on the train split of a dataset. We use ResNet-10[25] or WRN-28-10[75] as feature extractor backbone. This section will present the architecture and exact training procedure for our backbones.

**Network Architecture.** The architecture of our ResNet-10 and WRN-28-10 backbone is shown in Figure A1. Specifically, each convolutional layer is described as “ $n \times n$  conv,  $p$ ”, where  $n$  is the kernel size and  $p$  is the number of output channels. Convolutional layers with “/2” have a stride of 2 and are used to perform downsampling. The solid curved lines represent identity shortcuts, and the

dotted lines are projection shortcuts implemented by  $1 \times 1$  convolutions. The batch normalization and ReLU layers are omitted in Figure A1 to highlight the key structure of the two backbones.

**Pre-training Procedure.** The networks are trained from scratch with stochastic gradient descent in a fully-supervised manner, *i.e.*, minimizing cross-entropy loss on the train split of a dataset. Specifically the training is conducted on 90 epochs with early stopping using validation accuracy. We used batch size of 256 and image size of  $84 \times 84$ . For data augmentation, a random patch is sampled from an image, resized to  $84 \times 84$  and randomly flipped along horizontal axis before used for training. The initial learning rate is set to 0.1 and it is scaled down by factor of 10 every 30 epochs.

### A.5.2 Fine-Tuning

We consider linear, cosine and  $k$ -NN classifier for our fine-tuning experiments. In a  $K$ -way FSL problem, the detailed implementations for the classifier function  $f(\mathbf{x})$  are:

**Linear.**  $f(\mathbf{x}) = \mathbf{W}\mathbf{x} + \mathbf{b}$ , where  $\mathbf{x}$  is the input feature,  $\mathbf{W} \in \mathbb{R}^{K \times N}$  is the learnable weight parameter,  $N$  is the feature dimension and  $\mathbf{b} \in \mathbb{R}^K$  is the learnable bias parameter.

**Cosine.**  $f(\mathbf{x}) = \mathbf{W}\mathbf{x} / \|\mathbf{W}\| \|\mathbf{x}\|$ , where  $\mathbf{W} \in \mathbb{R}^{K \times N}$  is the learnable weight parameter. We implemented cosine classifier without using the bias term.

**$k$ -NN.** Our implementation of  $k$ -NN is similar to [57, 68]. For each of the  $K$  classes, we first calculated the average support set feature (centroid) denoted as  $\mathbf{x}_i, i \in \{1, \dots, K\}$ . The classifier output for class  $i$  is then given by  $f_i(\mathbf{x}) = -\|\mathbf{x} - \mathbf{x}_i\|^2$ . Notice that the prediction given by this classifier will be the nearest centroid.

We froze the backbone and used the average pooling layer output of  $\Omega$  to learn the classifier. The output logits from classifier functions are normalized using softmax to generate probability output  $P(y|\mathbf{x})$ . For linear and cosine classifier, we followed [13] and trained the classifier for 100 iteration with a batch size of 4. For fine-tuning baseline, we set the learning rate as  $1 \times 10^{-2}$  and weight decay as  $1 \times 10^{-3}$ . For IFSL, we set the learning rate as  $5 \times 10^{-3}$  and weight decay as  $1 \times 10^{-3}$ .  $k$ -NN classifier is non-parametric and can be initialized directly from support set.

### A.5.3 Meta-Learning

**MAML.** MAML [20] aims to learn an initialization of network parameters such that it can be fine-tuned within a few steps to solve a variety of few-shot classification tasks. When using pre-trained network with MAML, it has been shown that learning initialization of the backbone can lead to unsatisfactory performance [13, 58]. Therefore in our experiment, we froze the backbone and appended a 2-layer MLP with ReLU activation in between the hidden layers and a linear classifier after the average pooling layer of  $\Omega$ . The hidden dimension of the layers in MLP is the same as output dimension of  $\Omega$  (512 for ResNet-10 and 640 for WRN-28-10). The initialization of MLP and the linear classifier is meta-learnt using MAML. For hyper-parameters, we set the inner loop learning rate  $\alpha = 0.01$ , the outer loop learning rate  $\beta = 0.01$  and the number of adaptation steps as 20. For IFSL, we adopted the same hyper-parameter setting and set  $n=8$  for feature-wise and combined adjustment. Implementation-wise, we adopted the released code<sup>1</sup> from [13] and performed experiments on MAML without using first-order approximation. Following the implementation in [13], the model was trained on 10,000 randomly sampled tasks with model selection using validation accuracy. We used 2,000 randomly sampled tasks for validation and testing.

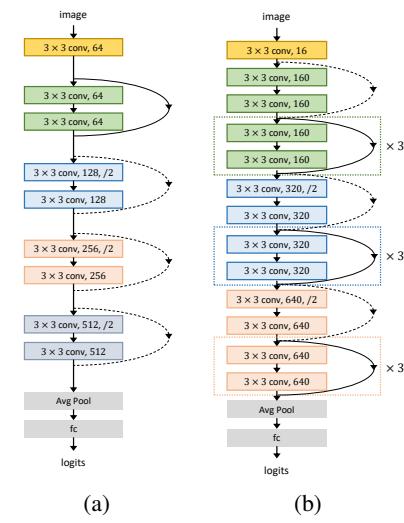


Figure A1: The architecture of our backbones: (a) ResNet-10 [25]; (b) WRN-28-10 [75].

<sup>1</sup><https://github.com/wyharveychen/CloserLookFewShot>

**MTL.** MTL [58] learns scaling and shifting parameters at each convolutional layer of the backbone. We used the MTL implementation released by the author<sup>2</sup> which adopts linear classifier. We integrated our ResNet-10 and WRN-28-10 backbones into the released code. The learning rate for scaling and shifting weights  $\phi_{SS}$  and initial classifier parameters was set to  $1 \times 10^{-4}$  uniformly. We set the inner loop learning rate for classifier as  $1 \times 10^{-2}$  and the inner loop update step as 100. For IFSL, we adopted the same hyper-parameter setting and set  $n=8$  for feature-wise and combined adjustment. We trained the MTL model on 10,000 randomly sampled tasks with model selection using validation accuracy and used 2,000 randomly sampled tasks for validation and testing. We used 3 RTX 2080 Ti for MTL experiments on WRN-28-10 backbone.

**LEO.** LEO [54] learns to generate classifier parameters conditioned on support set and the generated parameters are further fine-tuned within each FSL task. Our experiments were conducted on the released code of LEO<sup>3</sup> using linear classifier. Following author’s implementation, we saved the center cropped features from our pre-trained backbones and used the saved features to train LEO. For baseline, we used the hyper-parameter settings released by the author. For IFSL, we set  $n=8$  for feature-wise and combined adjustment and halved the outer loop learning rate compared to baseline. The model was trained up to 100,000 randomly sampled tasks from training split with early stopping using validation accuracy. We used 2,000 randomly sampled tasks for validation and testing.

**Matching Net.** Matching Net (MN) [66] is a metric-based method that learns a distance kernel function for  $k$ -NN. We used the Matching Net implementation in [13]. The implementation follows the setup in [66] and uses LSTM-based fully conditional embedding. We set the learning rate as 0.01 uniformly. For IFSL, we used  $n=16$  for feature-wise and combined adjustment. The model was trained using 10,000 randomly sampled tasks with model selection using validation accuracy. We used 2,000 randomly sampled tasks for validation and testing.

**SIB.** SIB [29] initializes classifier from support set and generates gradients conditioned on unlabelled query set features to update classifier parameters. We followed the SIB implementation released by the author<sup>4</sup> which uses cosine classifier. In the transductive setting, the query set size is set to 15. In the inductive setting, we used only 1 query sample randomly selected from the  $K$  classes in each episode. In terms of hyper-parameter settings, we took 3 synthetic gradient steps ( $K = 3$ ) for all our experiments. For baseline, the learning rate for SIB network and classifier was set to  $1 \times 10^{-3}$  following author’s implementation. For IFSL, we set the learning rate to  $5 \times 10^{-4}$  and used  $n=4$  for feature-wise and combined adjustment. In both transductive and inductive settings, we meta-trained SIB using 50,000 randomly sampled tasks with model selection using validation accuracy. We used 2,000 randomly sampled tasks for validation and testing.

## A.6 Additional Results

In this section, we include additional results on 1) **Conventional Acc** in Table A1 supplementary to Table 1; 2) **Hardness-Specific Acc** in Figure A2 for *miniImageNet* and Figure A3 for *tieredImageNet*, supplementary to Figure 5; 3) **CAM-Acc** in Table A2 supplementary to Figure 6; 4) **Cross-Domain Evaluation** in Table A3 supplementary to Table 3.

---

<sup>2</sup><https://github.com/yaoyao-liu/meta-transfer-learning>

<sup>3</sup><https://github.com/deepmind/leo>

<sup>4</sup>[https://github.com/hushell/sib\\_meta\\_learn](https://github.com/hushell/sib_meta_learn)

### A.6.1 Conventional Acc

Table A1: Supplementary to Table 1. Acc (%) and 95% confidence intervals averaged over 2000 5-way FSL tasks before and after applying three proposed implementations of adjustment. Specifically, “*ft*” refers to feature-wise adjustment, “*cl*” refers to class-wise adjustment and “*ft+cl*” refers to combined adjustment.

Method		ResNet-10				WRN-28-10				
		miniImageNet		tieredImageNet		miniImageNet		tieredImageNet		
		5-shot	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot	1-shot	
Fine-Tuning	Linear	76.38 ± 0.36	56.26 ± 0.47	81.01 ± 0.38	61.39 ± 0.47	79.79 ± 0.33	60.69 ± 0.45	85.37 ± 0.34	67.27 ± 0.49	
		<i>ft</i>	76.84 ± 0.36	57.37 ± 0.43	81.45 ± 0.38	61.88 ± 0.47	80.22 ± 0.31	60.84 ± 0.45	85.70 ± 0.33	67.94 ± 0.48
		<i>cl</i>	77.23 ± 0.34	59.45 ± 0.45	81.33 ± 0.38	62.60 ± 0.48	80.27 ± 0.32	62.15 ± 0.44	85.54 ± 0.33	68.11 ± 0.48
		<i>ft+cl</i>	77.97 ± 0.34	60.13 ± 0.45	82.08 ± 0.37	64.29 ± 0.48	80.97 ± 0.31	64.12 ± 0.44	86.19 ± 0.34	69.96 ± 0.46
Cosine		76.68 ± 0.36	56.40 ± 0.46	81.13 ± 0.39	62.08 ± 0.47	79.72 ± 0.33	60.83 ± 0.46	85.41 ± 0.34	67.30 ± 0.50	
		<i>ft</i>	76.83 ± 0.35	56.86 ± 0.44	81.34 ± 0.37	62.45 ± 0.47	79.80 ± 0.32	61.25 ± 0.44	85.74 ± 0.33	67.86 ± 0.46
		<i>cl</i>	76.99 ± 0.35	57.65 ± 0.45	81.42 ± 0.38	63.37 ± 0.48	79.96 ± 0.32	62.04 ± 0.45	85.77 ± 0.33	68.45 ± 0.46
		<i>ft+cl</i>	77.63 ± 0.34	59.84 ± 0.46	81.75 ± 0.38	64.47 ± 0.48	80.74 ± 0.32	63.76 ± 0.45	86.13 ± 0.33	69.36 ± 0.47
<i>k</i> -NN		76.63 ± 0.36	55.92 ± 0.46	80.85 ± 0.39	61.16 ± 0.48	79.60 ± 0.32	60.34 ± 0.45	84.67 ± 0.34	67.25 ± 0.52	
		<i>ft</i>	77.98 ± 0.34	60.71 ± 0.44	81.95 ± 0.36	65.66 ± 0.48	81.17 ± 0.31	64.87 ± 0.44	85.76 ± 0.34	71.00 ± 0.47
		<i>cl</i>	78.36 ± 0.35	61.32 ± 0.45	81.93 ± 0.37	65.71 ± 0.48	80.61 ± 0.31	64.43 ± 0.45	85.90 ± 0.33	70.08 ± 0.48
		<i>ft+cl</i>	78.42 ± 0.34	62.31 ± 0.44	81.98 ± 0.38	65.71 ± 0.47	81.08 ± 0.32	64.98 ± 0.43	86.06 ± 0.32	70.94 ± 0.49
MAML [20]		70.85 ± 0.38	56.59 ± 0.48	74.02 ± 0.41	59.17 ± 0.52	73.92 ± 0.36	58.02 ± 0.47	77.20 ± 0.38	61.40 ± 0.54	
		<i>ft</i>	73.84 ± 0.37	57.63 ± 0.47	80.19 ± 0.40	60.03 ± 0.51	78.82 ± 0.36	58.55 ± 0.48	84.74 ± 0.37	66.74 ± 0.52
		<i>cl</i>	73.01 ± 0.36	56.69 ± 0.48	78.41 ± 0.40	61.16 ± 0.53	76.22 ± 0.35	58.32 ± 0.46	81.74 ± 0.38	63.61 ± 0.51
		<i>ft+cl</i>	76.37 ± 0.37	59.36 ± 0.48	81.04 ± 0.39	63.88 ± 0.50	79.25 ± 0.34	62.84 ± 0.46	85.10 ± 0.39	67.70 ± 0.53
LEO [54]		74.49 ± 0.36	58.48 ± 0.48	80.25 ± 0.38	65.25 ± 0.51	75.86 ± 0.35	59.77 ± 0.47	82.15 ± 0.37	68.90 ± 0.49	
		<i>ft</i>	76.77 ± 0.35	60.52 ± 0.47	80.97 ± 0.36	65.44 ± 0.49	77.81 ± 0.34	61.81 ± 0.46	84.95 ± 0.36	69.59 ± 0.47
		<i>cl</i>	74.66 ± 0.36	58.62 ± 0.46	80.74 ± 0.37	65.37 ± 0.50	76.13 ± 0.35	60.22 ± 0.47	82.31 ± 0.37	69.23 ± 0.48
		<i>ft+cl</i>	71.91 ± 0.35	61.09 ± 0.47	81.43 ± 0.36	66.03 ± 0.48	77.72 ± 0.34	62.19 ± 0.45	85.04 ± 0.36	70.28 ± 0.47
MTL [58]		75.65 ± 0.35	58.49 ± 0.46	81.14 ± 0.36	64.29 ± 0.50	77.30 ± 0.34	62.99 ± 0.46	83.23 ± 0.37	70.08 ± 0.52	
		<i>ft</i>	77.17 ± 0.35	58.85 ± 0.44	82.01 ± 0.36	64.67 ± 0.47	79.40 ± 0.34	63.65 ± 0.45	84.76 ± 0.36	70.25 ± 0.49
		<i>cl</i>	77.10 ± 0.34	58.86 ± 0.45	82.34 ± 0.36	66.70 ± 0.51	79.29 ± 0.35	63.14 ± 0.46	86.21 ± 0.37	70.16 ± 0.50
		<i>ft+cl</i>	78.03 ± 0.33	61.17 ± 0.45	82.35 ± 0.35	65.72 ± 0.48	80.20 ± 0.33	64.40 ± 0.45	86.02 ± 0.35	71.45 ± 0.48
MN [66]		75.21 ± 0.35	61.05 ± 0.46	79.92 ± 0.37	66.01 ± 0.50	77.15 ± 0.36	63.45 ± 0.45	82.43 ± 0.37	70.38 ± 0.49	
		<i>ft</i>	75.52 ± 0.35	61.23 ± 0.45	80.18 ± 0.36	66.33 ± 0.49	77.80 ± 0.35	64.42 ± 0.46	83.82 ± 0.36	70.90 ± 0.50
		<i>cl</i>	75.40 ± 0.34	61.14 ± 0.44	80.04 ± 0.35	66.26 ± 0.50	77.23 ± 0.35	64.21 ± 0.47	82.77 ± 0.35	70.61 ± 0.51
		<i>ft+cl</i>	76.73 ± 0.34	62.64 ± 0.46	80.79 ± 0.35	67.30 ± 0.48	78.55 ± 0.36	64.89 ± 0.44	84.03 ± 0.36	71.41 ± 0.49
SIB [29] (transductive)		78.88 ± 0.35	67.10 ± 0.56	85.09 ± 0.35	77.64 ± 0.58	81.73 ± 0.34	71.31 ± 0.56	88.19 ± 0.34	81.97 ± 0.56	
		<i>ft</i>	79.58 ± 0.35	67.94 ± 0.55	85.12 ± 0.35	77.68 ± 0.57	82.00 ± 0.34	71.95 ± 0.56	88.20 ± 0.34	82.01 ± 0.56
		<i>cl</i>	79.04 ± 0.33	67.77 ± 0.55	85.22 ± 0.35	77.72 ± 0.56	81.93 ± 0.35	71.66 ± 0.56	88.21 ± 0.33	82.01 ± 0.54
		<i>ft+cl</i>	80.32 ± 0.35	68.85 ± 0.56	85.43 ± 0.35	78.03 ± 0.57	83.21 ± 0.33	73.51 ± 0.56	88.69 ± 0.33	83.07 ± 0.52
SIB [29] (inductive)		75.64 ± 0.36	57.20 ± 0.57	81.69 ± 0.34	65.51 ± 0.56	78.17 ± 0.35	60.12 ± 0.56	84.96 ± 0.36	69.20 ± 0.58	
		<i>ft</i>	76.23 ± 0.35	58.67 ± 0.56	82.04 ± 0.35	66.69 ± 0.57	79.34 ± 0.35	61.77 ± 0.56	85.24 ± 0.36	70.05 ± 0.57
		<i>cl</i>	76.61 ± 0.35	58.12 ± 0.55	82.21 ± 0.35	66.28 ± 0.56	79.11 ± 0.35	61.25 ± 0.55	85.63 ± 0.34	69.90 ± 0.57
		<i>ft+cl</i>	77.68 ± 0.34	60.33 ± 0.54	82.75 ± 0.35	67.34 ± 0.55	80.05 ± 0.34	63.14 ± 0.54	86.14 ± 0.34	71.45 ± 0.55

### A.6.2 Hardness-Specific Acc

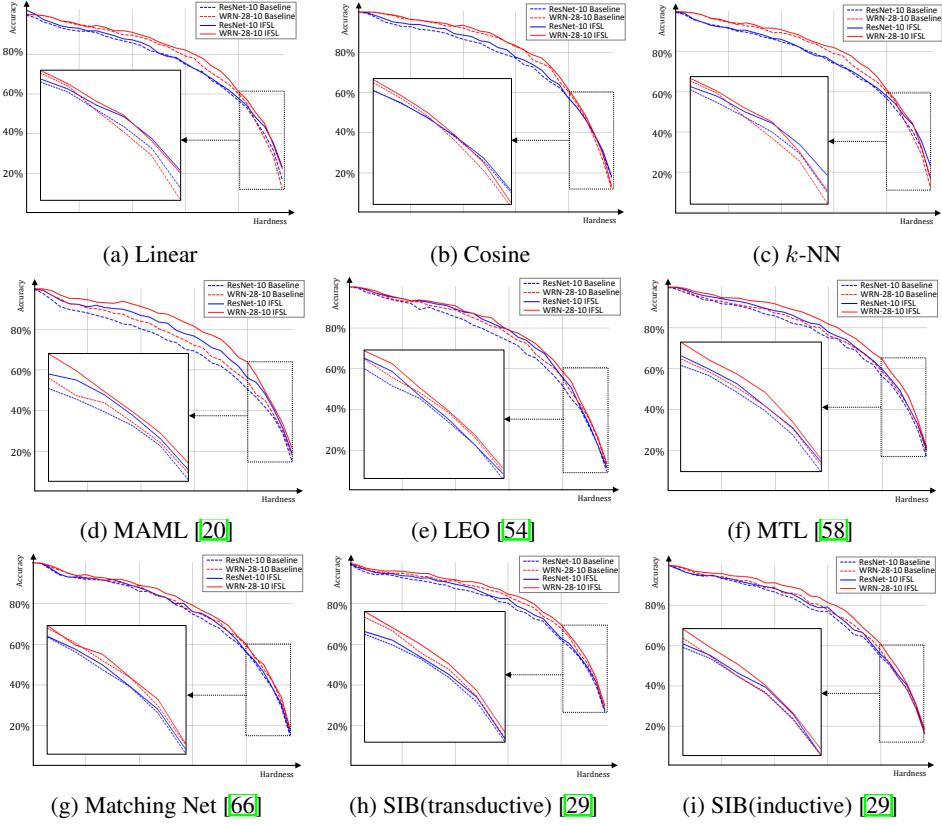


Figure A2: Supplementary to Figure 5. Hardness-specific Acc of 5-shot fine-tuning and meta-learning on *miniImageNet*.

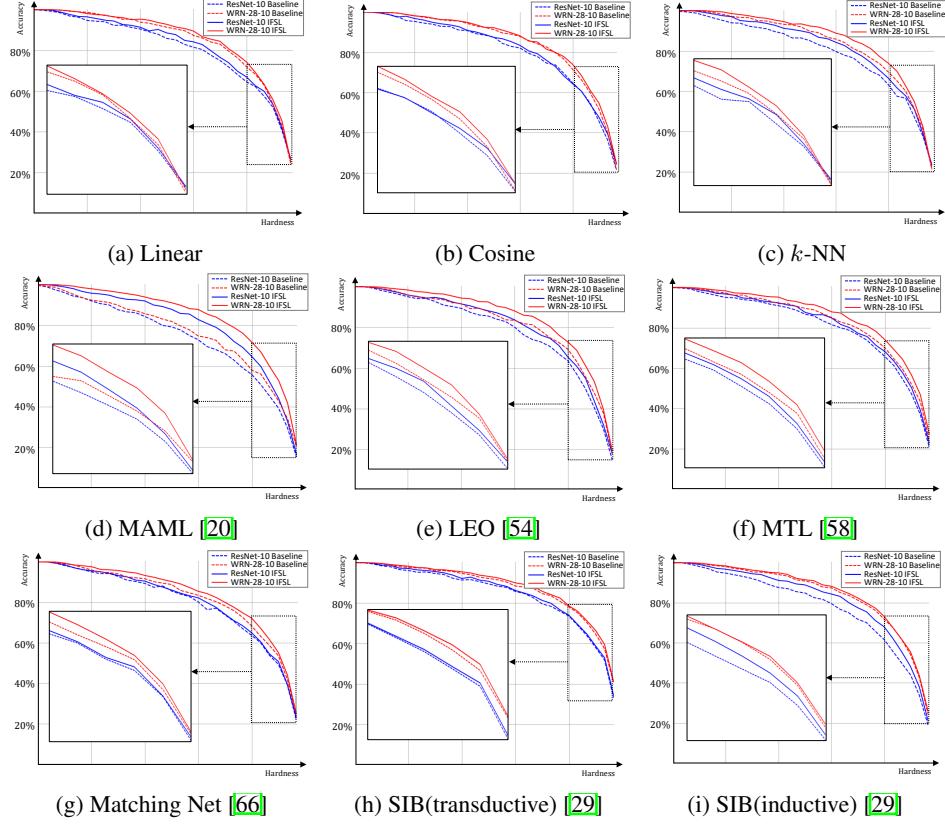


Figure A3: Supplementary to Figure 5. Hardness-specific Acc of 5-shot fine-tuning and meta-learning on *tieredImageNet*.

### A.6.3 CAM-Acc

Table A2: Supplementary to Figure 6. CAM-Acc (%) on fine-tuning and meta-learning. We used combined adjustment for IFSL.

Method	ResNet-10				WRN-28-10				
	miniImageNet		tieredImageNet		miniImageNet		tieredImageNet		
	5-shot	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot	1-shot	
Fine-Tuning	Linear	29.02 ± 0.38	25.22 ± 0.38	31.62 ± 0.38	31.05 ± 0.39	25.99 ± 0.35	24.74 ± 0.34	30.17 ± 0.36	29.76 ± 0.37
	+IFSL	29.85 ± 0.37	26.67 ± 0.38	31.75 ± 0.38	31.43 ± 0.37	26.02 ± 0.37	24.96 ± 0.36	32.57 ± 0.36	30.64 ± 0.38
	Cosine	28.10 ± 0.37	27.12 ± 0.38	29.82 ± 0.37	28.54 ± 0.38	27.54 ± 0.38	25.73 ± 0.37	32.60 ± 0.35	31.21 ± 0.36
	+IFSL	28.18 ± 0.37	27.26 ± 0.38	31.38 ± 0.37	28.70 ± 0.40	27.82 ± 0.38	25.85 ± 0.38	33.65 ± 0.37	31.66 ± 0.35
	<i>k</i> -NN	27.96 ± 0.37	26.65 ± 0.37	32.25 ± 0.39	30.36 ± 0.39	24.15 ± 0.34	23.30 ± 0.33	23.91 ± 0.34	21.99 ± 0.33
	+IFSL	28.15 ± 0.37	26.81 ± 0.37	32.75 ± 0.39	30.84 ± 0.39	25.23 ± 0.36	24.14 ± 0.35	28.04 ± 0.37	26.46 ± 0.37
Meta-Learning	MAML [20]	29.43 ± 0.37	27.39 ± 0.38	32.72 ± 0.40	32.14 ± 0.40	27.56 ± 0.36	26.46 ± 0.36	34.39 ± 0.40	31.07 ± 0.39
	+IFSL	30.06 ± 0.38	28.42 ± 0.38	32.93 ± 0.40	32.24 ± 0.39	27.61 ± 0.36	26.91 ± 0.38	34.57 ± 0.41	31.22 ± 0.40
	LEO [54]	30.24 ± 0.38	28.56 ± 0.37	31.64 ± 0.38	29.88 ± 0.37	29.15 ± 0.38	27.86 ± 0.38	31.27 ± 0.37	29.73 ± 0.38
	+IFSL	30.67 ± 0.37	28.76 ± 0.37	32.01 ± 0.38	30.65 ± 0.37	29.20 ± 0.37	28.45 ± 0.38	31.98 ± 0.39	30.32 ± 0.38
	MTL [58]	31.45 ± 0.39	30.13 ± 0.39	33.52 ± 0.39	33.11 ± 0.39	30.56 ± 0.39	29.78 ± 0.40	33.13 ± 0.39	32.35 ± 0.39
	+IFSL	34.21 ± 0.39	31.59 ± 0.40	33.67 ± 0.38	33.50 ± 0.39	31.78 ± 0.39	30.12 ± 0.39	33.30 ± 0.39	32.64 ± 0.39
	MN [66]	28.50 ± 0.38	28.42 ± 0.39	32.55 ± 0.40	31.88 ± 0.39	24.93 ± 0.38	25.34 ± 0.39	34.87 ± 0.37	29.10 ± 0.38
SIB [29]	+IFSL	28.68 ± 0.38	28.77 ± 0.38	32.67 ± 0.40	32.10 ± 0.40	27.93 ± 0.37	25.81 ± 0.37	35.47 ± 0.41	30.71 ± 0.39
	(transductive)	32.10 ± 0.39	31.19 ± 0.39	32.16 ± 0.39	30.49 ± 0.39	28.32 ± 0.37	26.76 ± 0.38	31.02 ± 0.36	28.43 ± 0.38
	+IFSL	32.14 ± 0.39	31.34 ± 0.39	34.31 ± 0.40	32.59 ± 0.40	31.54 ± 0.38	29.82 ± 0.36	32.33 ± 0.37	30.26 ± 0.38
SIB [29]	(inductive)	31.26 ± 0.38	30.56 ± 0.39	31.35 ± 0.40	30.48 ± 0.39	29.76 ± 0.38	28.02 ± 0.37	29.45 ± 0.39	27.98 ± 0.39
	+IFSL	31.46 ± 0.39	30.78 ± 0.40	31.56 ± 0.39	30.89 ± 0.40	30.23 ± 0.37	28.75 ± 0.39	30.07 ± 0.40	28.57 ± 0.39

#### A.6.4 Cross-Domain Evaluation

Table A3: Supplementary to Table 3. Acc (%) and 95% confidence interval averaged over 2000 5-way FSL tasks on cross-domain evaluation. Specifically, “*ft*” refers to feature-wise adjustment, “*cl*” refers to class-wise adjustment and “*ft+cl*” refers to combined adjustment.

Method		ResNet-10		WRN-28-10	
		5-shot	1-shot	5-shot	1-shot
Fine-Tuning	Linear	58.84 ± 0.41	42.25 ± 0.42	62.12 ± 0.40	42.89 ± 0.41
		60.12 ± 0.39	42.30 ± 0.41	63.13 ± 0.39	43.39 ± 0.40
		60.51 ± 0.40	42.43 ± 0.42	62.95 ± 0.39	44.21 ± 0.40
	Cosine	60.65 ± 0.39	45.14 ± 0.40	64.15 ± 0.38	45.64 ± 0.39
		58.30 ± 0.39	40.47 ± 0.40	60.21 ± 0.39	42.12 ± 0.39
		58.32 ± 0.39	41.01 ± 0.40	61.16 ± 0.38	42.35 ± 0.41
	<i>k</i> -NN	58.68 ± 0.39	40.67 ± 0.41	61.87 ± 0.40	43.23 ± 0.40
		60.23 ± 0.38	42.78 ± 0.40	62.49 ± 0.38	45.12 ± 0.39
		57.18 ± 0.40	38.44 ± 0.37	59.31 ± 0.41	40.53 ± 0.42
Meta-Learning	MAML [20]	59.44 ± 0.39	43.49 ± 0.40	62.48 ± 0.39	45.68 ± 0.43
		58.37 ± 0.39	43.20 ± 0.41	62.04 ± 0.39	45.36 ± 0.40
		59.59 ± 0.40	43.45 ± 0.40	62.45 ± 0.40	45.72 ± 0.40
	LEO [54]	51.09 ± 0.43	37.20 ± 0.46	55.04 ± 0.42	39.06 ± 0.47
		54.95 ± 0.44	37.34 ± 0.47	59.57 ± 0.44	39.25 ± 0.46
		53.62 ± 0.43	38.13 ± 0.47	56.80 ± 0.45	40.32 ± 0.48
	MTL [58]	56.71 ± 0.46	40.36 ± 0.46	60.89 ± 0.45	42.16 ± 0.47
		56.52 ± 0.46	39.21 ± 0.53	56.66 ± 0.48	41.45 ± 0.54
		56.77 ± 0.48	39.72 ± 0.54	62.95 ± 0.47	45.46 ± 0.55
	MN [66]	56.73 ± 0.47	40.12 ± 0.55	56.90 ± 0.47	41.93 ± 0.56
		61.27 ± 0.46	42.79 ± 0.52	63.30 ± 0.47	43.81 ± 0.56
		56.61 ± 0.42	41.56 ± 0.43	56.89 ± 0.41	43.15 ± 0.44
(transductive)	SIB [29]	61.34 ± 0.41	42.90 ± 0.43	63.49 ± 0.40	45.28 ± 0.44
		60.62 ± 0.41	42.87 ± 0.42	62.94 ± 0.40	45.57 ± 0.43
		62.39 ± 0.40	44.51 ± 0.43	65.00 ± 0.40	46.67 ± 0.43
	SIB [29]	53.39 ± 0.46	40.34 ± 0.56	53.08 ± 0.45	42.04 ± 0.57
		54.22 ± 0.46	40.62 ± 0.57	54.97 ± 0.47	42.52 ± 0.58
		53.72 ± 0.47	40.42 ± 0.56	53.43 ± 0.45	42.19 ± 0.56
	(inductive)	56.03 ± 0.45	41.68 ± 0.54	58.69 ± 0.44	43.58 ± 0.56
		60.60 ± 0.46	45.87 ± 0.55	62.60 ± 0.49	49.16 ± 0.58
		61.12 ± 0.45	46.64 ± 0.55	63.15 ± 0.47	49.78 ± 0.56
		60.70 ± 0.46	46.14 ± 0.56	63.02 ± 0.48	49.43 ± 0.57
		62.07 ± 0.44	47.07 ± 0.53	64.07 ± 0.49	50.71 ± 0.54
	(inductive)	59.06 ± 0.42	41.48 ± 0.43	59.94 ± 0.42	43.27 ± 0.44
		59.45 ± 0.41	41.98 ± 0.44	60.33 ± 0.44	43.61 ± 0.45
		59.32 ± 0.42	41.67 ± 0.43	60.46 ± 0.43	43.52 ± 0.45
		59.89 ± 0.41	43.20 ± 0.43	61.45 ± 0.43	44.27 ± 0.44