EARLY LAYER READOUTS FOR ROBUST KNOWLEDGE DISTILLATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Domain generalization (DG) aims to learn a model that can generalize to unseen i.e. out-of-distribution (OOD) test domain. While large-capacity networks trained with sophisticated DG algorithms tend to achieve high robustness, they tend to be impractical in deployment. Typically, Knowledge distillation (KD) can alleviate this via an efficient transfer of knowledge from a robust teacher to a smaller student network. Throughout our experiments, we find that vanilla KD already provides strong OOD performance, often outperforming standalone DG algorithms. Motivated by this observation, we propose an adaptive distillation strategy that utilizes early layer predictions and uncertainty measures to learn a meta network that effectively rebalances supervised and distillation losses as per sample difficulty. Our method adds no inference overhead and consistently outperforms canonical ERM, vanilla KD, and competing DG algorithms across OOD generalization benchmarks.

1 Introduction

Deploying machine learning models in real-world scenarios requires robustness to distribution shifts (Koh et al., 2021; Huang et al., 2021), often referred to as domain generalization (DG) (Zhao et al., 2020; Robey et al., 2021) or out-of-distribution (OOD) generalization (Wald et al., 2021; Montasser et al., 2024). While high-capacity models trained with specialized DG algorithms (Gulrajani & Lopez-Paz, 2020) can achieve strong robustness, they are often prohibitively expensive in terms of computation and memory, making them impractical for deployment in resource-constrained environments.

Knowledge distillation (KD) (Hinton et al., 2015a;b; Lopes et al., 2017), has emerged as a standard approach for improving efficiency by transferring knowledge from a large teacher model to a compact student model. Beyond efficiency aspect of KD, it has recently been explored for improving OOD robustness, where vanilla KD, as shown in (Zhou et al., 2022; 2023; Huang et al., 2023) tends to yield better OOD performance compared to models trained solely with DG algorithms. However, there is still room for improvement, as KD typically treates all samples uniformly and often overrelies on the teacher's dark knowledge, making it prone to teacher-specific biases. Moreover, exisiting works that adapt KD for domain generalization primarily focus on using adversarially trained teachers (Nasery et al., 2022), multimodal teacher networks for additional supervision Huang et al. (2023) or ensemble of domain-specific teachers Zhao et al. (2025), leaving the role of the student network underexplored.

To address these limitations, we propose an adaptive distillation framework that modulates distillation loss based on sample difficulty via with a lightweight *forecaster* meta-network. The forecaster leverages early layer representations and uncertainty measures to estimate sample difficulty and dynamically reweight the distillation loss. This enables the student to selectively trust the teacher where appropriate while emphasizing supervision from ground-truth labels for harder or biased samples. Crucially, our design introduces no additional inference-time overhead as the forecaster is discarded post-training, making it well-suited for practical deployment.

Our work makes the following contributions:

We identify the limitations, and opportunities for improvement in standard KD under domain shifts, noting that uniform treatment of samples and blind reliance on the teacher hinder OOD robustness.

- We propose an adaptive distillation framework with a *forecaster* meta-network that leverages early readouts to dynamically assign instance-specific weights in the loss function based on sample difficulty.
- We show that our approach improves student robustness with affecting deployment efficiency: it adds no inference-time overhead while consitently improving OOD generalization across multiple benchmarks.

2 BACKGROUND AND RELATED WORK

Instance-Specific Learning. Prior works have explore instance-specific learning to improve neural network training under noisy conditions, via instance-specific parameters such as temperature, smoothing factors, or weights. Saxena et al. (2019); Wang et al. (2018); Algan & Ulusoy (2021) introduce learnable sample and class weights to control the importance of each sample in the learning process depending on sample reliability and label noise. Ren et al. (2019); Shu et al. (2019); Raghu et al. (2021); Jain et al. (2024) adopt meta-learning (Finn et al., 2017) using small unbiased meta-samples to learn weighing functions to obtain instance-specific weights to address class imbalance, label noise and robustness. In the context of knowledge distillation, Zhao et al. (2021) propose a curriculum-based distillation with instance-level sequence learning. Iliopoulos et al. (2022) presents reweighing strategy for the student loss in knowledge distillation with unlabeled data, to eliminate potential biases from the teacher network. Sivasubramanian et al. (2022) present a bi-level objective to learn instance-specific combination of teacher-matching and supervised objectives to learn student models that are more accurate. In this work, we introduce a meta-network which guides the student via instance-specific weighing in the KD objective. Further, we interleave the training of the meta-network with the student, without requiring complex meta-learning.

Early Readouts. Prior works on Early Readouts majorly focus on early-exiting (Han et al., 2021; Xu & McAuley, 2023; Laskaridis et al., 2021; Matsubara et al., 2022) with the aim to reduce inference cost, allowing samples to "exit" at intermediate layers via auxiliary classifiers. These include dynamic early-exiting and static early-exiting mechanisms. Dynamic methods focus on balancing trade-off between speed and accuracy at inference, by early-exit mechanisms based on dynamics of internal classifiers, such as calibrated prediction confidence and entropy(Xin et al., 2020; Liu et al., 2020; Schwartz et al., 2020; Zhou et al., 2020), class mean of sample predictions (Görmez et al., 2022). In contrast to dynamic mechanisms, Sun et al. (2022) propose a hash-based early exiting for sequence learning tasks, where tokens are assigned to fixed exiting layers using a hash function. In the context of KD, Tiwari et al. (2023) use early readout errors to detect spurious feature reliance, and propose a weighing scheme to reweigh the distillation loss to reduce feature-specific bias. In this work, we re-purpose early readouts not for exiting, but as signals to guide our meta-network, forecaster, in assigning instance-specific weights in KD for OOD robustness, while avoiding the need for handcrafted weighing functions.

Distillation-based Domain Generalization. Distillation has shown promise in OOD generalization by allowing knowledge transfer from a robust teacher network, as opposed to training student network solely on a DG Algorithm (Wang et al., 2021; Huang et al., 2023). However, prior works on KD for Domain Generalization focus on teacher network or the teacher-student interaction. Wang et al. (2021) propose gradient based regularization to lower the mapping difficulty from the teacher to the student. Nasery et al. (2022) utilize adversarially fine-tuned teacher networks to improve knowledge transfer to student for OOD generalization. Huang et al. (2023) leverage CLIP teacher model along with a proposed text-based regularization scheme to enable better transfer from teacher. Zhao et al. (2025) leverage domain-specific teachers to improve student generalizability in an online KD setting. *In contrast to prior works, our method explores student-centric adaptation, leveraging early layer prediction confidences to navigate the distillation process*.

3 NOTATION AND PROBLEM SETUP

Notations. The first set of n natural numbers $\{1, 2, \ldots, n\}$ is denoted by [n]. The n-dimensional real vector space is denoted by \mathbb{R}^n . Vectors are typeset in lowercase bold (e.g., \mathbf{x}); matrices are in uppercase bold (e.g., \mathbf{X}); and elements are referenced by subscripts (e.g., \mathbf{x}_i , \mathbf{X}_{ij}). When needed for clarity, elements will be referenced by subscripts on square brackets (e.g., $[\mathbf{x}_1]_i$, $[\mathbf{X}_2]_{ij}$). We denote the sigmoid function by $\sigma(x) = (1 + e^{-x})^{-1}$.

Knowledge Distillation. Knowledge distillation (Hinton et al., 2015b) provides an efficient framework for training compact models without the need to optimize over large-scale networks, thereby reducing both memory footprint and computational overhead. In this paradigm, the compact model (student) is trained to align with both the ground-truth labels and the predictive behavior of a larger reference model (teacher). Formally, the training objective for the student consists of two components:

- 1. **Supervised Loss** \mathcal{L}_{CE} , which measures the discrepancy between the student's predictions and the ground-truth labels using cross-entropy.
- 2. **Distillation Loss** \mathcal{L}_{KD} , defined as the Kullback-Leibler divergence between the student's and teacher's output distributions. This term encourages the student to replicate the softened predictive probabilities of the teacher, which encode richer information than hard labels alone and thus facilitate more effective knowledge transfer.

Problem Setting. Consider a standard supervised multi-classification setting with inputs $x \in \mathcal{X}$, outputs $y \in \mathcal{Y}$, and the training data $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ where $y_i \in [K]$, K denoting the total number of classes and $\mathbf{x}_i \in \mathbb{R}^d$ denoting the i-th input feature. Suppose, we train a student network $\mathcal{M}_{\boldsymbol{\theta}^{(S)}}$ parameterized by $\boldsymbol{\theta}^{(S)}$ and let the corresponding teacher network consists of L layers, where the output representation for the ℓ -th layer $\ell \in L$ is given by $\mathbf{z}_{\boldsymbol{\theta}^{(S)}}^\ell = h_{\boldsymbol{\theta}^{(S)}}^\ell(\mathbf{x}_i) \in \mathbb{R}^{d_\ell}$ where d_ℓ denotes the ℓ -th layer dimension. The final layer output logits from student network corresponding to an input \mathbf{x}_i is denoted by $\mathbf{z}_{\boldsymbol{\theta}^{(S)}}^\ell := \mathcal{M}_{\boldsymbol{\theta}^{(S)}}(\mathbf{x}_i)$.

Loss Formulation. In KD, the student network $\mathcal{M}_{\boldsymbol{\theta}^{(S)}}$ is trained by minimizing the standard cross entropy loss between the model predictions and ground truth over the training data \mathcal{D} as per (1), while the teacher predictions are used to minimize the knowledge distillation term (2), defined as the KL divergence between teacher and student logits $\mathcal{M}_{\boldsymbol{\theta}^{(T)}}(\mathbf{x}_i)$ and $\mathcal{M}_{\boldsymbol{\theta}^{(S)}}(\mathbf{x}_i)$.

$$\mathcal{L}_{CE}^{(S)} = -\frac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{K} \mathbf{1}_{(y_i = j)} \log p_{\boldsymbol{\theta}^{(S)}}(\mathbf{x}_i)_j \quad (1) \quad \mathcal{L}_{KD} = \frac{1}{N} \sum_{i=1}^{N} \mathbb{D}_{KL}(p_{\boldsymbol{\theta}^{(T)}}(\mathbf{x}_i) \parallel p_{\boldsymbol{\theta}^{(S)}}(\mathbf{x}_i)), \quad (2)$$

where $p_{\theta}(\mathbf{x}_i)_j = \frac{\exp(\mathcal{M}_{\theta}(\mathbf{x}_i)_j/\tau)}{\sum_{k=1}^K \exp(\mathcal{M}_{\theta}(\mathbf{x}_i)_k/\tau)}$ denotes j-th class probability where $j \in [K]$. Here, τ denotes the temperature used to soften the output probabilities of both the student and the teacher. Generally, a higher temperature allows for smoother distributions, which capture richer inter-class relationship from the teacher (Hinton et al., 2015a). This allows the student to capture more nuanced relationships between classes rather than focusing on a single class representing the highest probability.

Student Training Objective. The student network's loss is denoted by $\mathcal{L}_{\mathrm{tot}}^{(S)}(\cdot,\cdot;\boldsymbol{\theta}^{(S)})$, where

$$\mathcal{L}_{\text{tot}}^{(S)}(\boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{Y}}; \boldsymbol{\theta}^{(S)}) = \alpha \cdot \mathcal{L}_{CE}^{(S)} + \beta \cdot \mathcal{L}_{KD}$$
(3)

Here, $\mathcal{L}_{\text{tot}}^{(S)}(\cdot,\cdot;\boldsymbol{\theta}^{(S)})$ is expressed as a convex combination of $\mathcal{L}_{CE}^{(S)}$ and \mathcal{L}_{KD} , with $\alpha \in \mathbb{R}$ and $\beta \in \mathbb{R}$ being seen as two degrees of freedom, controlling the contribution of each loss component (Srivastava et al., 2015). The two degrees of freedom can be reduced to a single degree of freedom by bounding α and β with the condition: $\alpha + \beta = 1$.

Moreover, in the absence of labeled data to train the student model, the distillation reduces to a soft-distillation setting, where $\alpha=0$. In this case, the soften teacher probabilities form the only source of supervision for the student model (Lopes et al., 2017).

Student capacity limits and Teacher Confidence. The student's ability to learn the teacher's dark knowledge is inherently bounded, and the student may fail to capture richer information beyond a threshold due to limited capacity. On the other hand, larger teachers are over-confident, yielding higher target logits and lower variance in predictions, resulting in less distinctive incorrect-class softmax probabilities. In this case, if distillation temperature is increased, teacher guidance becomes weaker, smoothing strengthens due to the softened distribution, and class discriminability measured by variance of incorrect-class probability initially rises then falls off. Together, both the student capacity and teacher confidence restrict the effectiveness of knowledge transfer (Li et al., 2022).

3.1 OOD GENERALIZATION IN DISTILLED MODELS

Bottlenecks with Vanilla KD. While distillation serves as a good approach for model compression, the student model can be interpreted as a small clone of the teacher over-fitted to the teacher's learned patterns. The over-fitting and the student model's limited capacity typically fails to generalize well with OOD samples at inference time (Yue et al., 2023). While there exist OOD generalization algorithms targeted towards OOD robustness in neural network (Gulrajani & Lopez-Paz, 2020), the lack of representation capacity affects the ability of the student to directly benefit from these algorithms.

A general solution to improve generalization ability of the student is to train the teacher network with domain generalization (Gulrajani & Lopez-Paz, 2020) algorithms or adversarial training (?) (Nasery et al., 2022) for OOD robustness. During distillation, the robust teacher can act as a "shortcut" that helps the student model bypass the requirement of complex, robust training. The student model inherits the teacher's ability to produce well-calibrated outputs for seen and unseen input instances. This helps the student to be robust to distribution shifts in a deployment environment.

3.2 EARLY READOUTS

Early-layer confidences in a neural network offer valuable insights about an input sample (Baldock et al., 2021). Low confidences at these layers can suggest that the sample is complex or ambiguous, potentially making it harder for the model to classify accurately. Conversely, higher early-layer confidences indicate that the low-level representations align more closely with a particular class, increasing the likelihood that the student network will classify it correctly. By providing additional supervision, such as hard labels, for ambiguous samples, the model can improve its learning of overlapping or confusing features across classes, ultimately enhancing overall performance. (Tiwari et al., 2024) utilize early exit information from neural networks.

4 Proposed Methodology

Our overall architecture consists of the teacher-student setup, with alterations to the student network. We incorporate auxillary networks (internal classifiers) on intermediate layers' output representations, say denoted as $\mathcal{E} \subseteq [L]$, where L denotes the number of layers in the model.

Figure 1 provides a brief overview of our approach where we incorporate early-layer predictions and uncertainty measures from the student model to dynamically weigh the individual loss components of the distillation objective. In the sections to follow, we elaborate on the auxiliary networks and the *forecaster* in detail.

4.1 AUXILIARY NETWORK: EARLY LAYER CONFIDENCE

Design of Auxiliary Networks. For each early layer $\ell \in \mathcal{E}$, we instantiate an auxiliary predictor $\mathbf{A}_{\mathrm{aux}}^{(\ell)}(\bullet; \varphi_{\ell}): \mathbb{R}^{d_{\ell}} \to \Delta^K$, with the ℓ -th auxiliary classifier parameterized by φ_{ℓ} . The input to $\mathbf{A}_{\mathrm{aux}}^{(\ell)}$ is the ℓ -th layer's output representation $\mathbf{z}_{\varrho(S)}^{\ell}$. The role of the auxiliary classifier is to encode how well the truncated feature stack up to layer ℓ differentiates among the label space \mathcal{Y} , i.e., the extent to which intermediate layers capture discriminative information over class labels. Such auxiliary predictors have been shown to be effective both for improving gradient flow and feature discriminability during training (Szegedy et al., 2015; Lee et al., 2015), as well as for posthoc analysis of representation quality via probing (Alain & Bengio, 2016).

Training Objective of Auxiliary Networks. On the training split \mathcal{D} , the objective of each auxiliary encoder is to minimize the standard classification loss:

$$\mathcal{L}_{\text{aux}}^{(\ell)}(\mathbf{z}_{\boldsymbol{\theta}^{(S)}}^{\ell}; \boldsymbol{\varphi}_{\ell}) = -\sum_{i=1}^{N} \sum_{j=1}^{K} \mathbf{1}_{\{y_i = j\}} \log \frac{\exp(\mathbf{A}_{\text{aux}}^{(\ell)}(\mathbf{z}_{\boldsymbol{\theta}^{(S)}}^{\ell}; \boldsymbol{\varphi}_{\ell})_j)}{\sum_{k=1}^{K} \exp(\mathbf{A}_{\text{aux}}^{(\ell)}(\mathbf{z}_{\boldsymbol{\theta}^{(S)}}^{\ell}; \boldsymbol{\varphi}_{\ell})_k)}.$$
 (4)

4.2 Forecaster Meta-Network

Deep neural networks are prone to overfitting under label noise or class imbalance, where fixed re-weighting schemes often fail due to their reliance on handcrafted weighting functions and hyperparameters. To address this, we introduce a *forecaster* module that adaptively maps sample-level statistics to re-weigh the influence of the teacher network in the distillation process. Parameterized

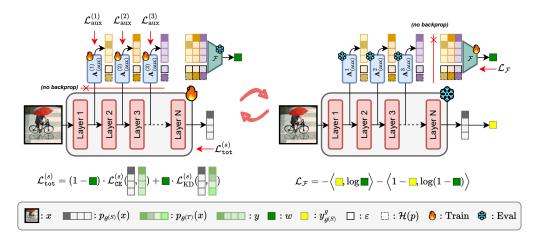


Figure 1: Overview of adaptive knowledge distillation that leverages early layer readout predictions and uncertainty measures to learn a meta-network *forecaster* which modulates contribution of ground-truth supervision and teacher supervision in the distillation process at an instance level. Here, the student network, early layer auxiliary networks and *forecaster* are trained in an interleaved fashion. First, the student backbone and the auxiliary network are trained on the classification objective, while the *forecaster* remains frozen, for a fixed number of train minibatches. Next, for a fixed number meta-validation set from the training domains is used to train the forecaster with the objective to estimate student correctness. At test time, the *forecaster* and auxiliary networks are discarded, resulting in our method requiring the same computational resources as a vanilla KD during inference.

as a lightweight neural network and optimized jointly with the model, the forecaster provides a flexible data-driven mechanism that generalizes beyond handcrafted weighting rules. We denote the forecaster as $\mathcal{F}(\bullet; \psi^f)$ defined as a meta-network parameterized by ψ^f .

Confidence Margin. The confidence margin of the auxiliary networks denoted by ε acts a good signal of uncertainty in the predictions at the early layers (Tiwari et al., 2024; Xin et al., 2020; Liu et al., 2020) where $p_{\max}(\varphi_l) = \max(p_1(\varphi_l), p_2(\varphi_l), \ldots, p_K(\varphi_l))$ and $p_K(\cdot)$ indicates the probability of the auxiliary layer ℓ 's output being classified as label K.

$$\varepsilon^{(\ell)} = p_{\text{max}}(\cdot) - \max(p_i(\cdot) \mid i \neq \arg\max(p_i(\cdot)))$$
 (5)

Particularly, *forecaster* is fed logits (z), Entropy $(\mathcal{H}(p))$ and Confidence Margin (ε) of the auxiliary network predictions as uncertainty indicators of auxiliary network, highlighting prediction confidence and randomness.

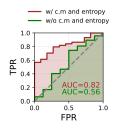


Figure 2: ROC curves showcasing the impact of incorporating confidence margin (ε) and entropy $\mathcal{H}(p)$ features from early layers into the forecaster.

Training Forecaster $\mathcal{F}(ullet;\psi^f)$:The Forecaster is trained in

conjunction to the student $\mathcal{M}_{\theta^{(S)}}$ and the auxiliary networks $\mathbf{A}_{\mathsf{aux}}^{(\ell)}(\bullet;\varphi)$ using a binary classification objective. The binary classification task for the forecaster is to determine whether student model correctly predicts an input instance. Higher output probability from the forecaster indicates student is more likely to correctly classify the sample, indicating an easier sample. We utilize the output probability of the forecaster directly to weigh the KL divergence between the student and the teacher, forcing the student to mimic teacher for such samples and focus more on ground-truth supervision when the forecaster output probability is smaller. To train the forecaster, we minimize $\mathcal{L}_{\text{focs}}(\theta^{(S)};\psi^f)$ on the validation split $\mathcal{D}_v\subseteq\mathcal{D}$

$$\mathcal{L}(\boldsymbol{\theta}^{(S)}, \boldsymbol{\psi}_f) = -\left\langle \mathbf{y}_{\boldsymbol{\theta}^{(S)}}^g, \log \mathbf{w}_{\boldsymbol{\psi}_f} \right\rangle - \left\langle \mathbf{1} - \mathbf{y}_{\boldsymbol{\theta}^{(S)}}^g, \log(\mathbf{1} - \mathbf{w}_{\boldsymbol{\psi}_f}) \right\rangle, \tag{6}$$

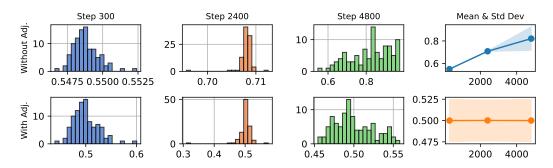


Figure 3: Distribution of forecaster outputs at different training steps, comparing the unadjusted case (**top row**) with post-hoc adjustment (**bottom row**). The rightmost plots summarize the deviations across epochs. Post-hoc adjustment stabilizes the output distribution, negating drift towards skewed values.

where $\mathbf{y}_{\boldsymbol{\theta}^{(S)}}^g \in \{0,1\}^N$ is the binary vector of ground-truth labels for the student model predictions, with entries $y_{\boldsymbol{\theta}^{(S)},i}^g = \mathbf{1}_{(\hat{y}_i = y_i)}$ for each instance $i \in [N]$, and $\mathbf{w}_{\boldsymbol{\psi}_f} \in (0,1)^N$ is the vector of forecaster outputs parameterized by $\boldsymbol{\psi}_f$.

Design of Forecaster. The *forecaster* is implemented as a lightweight neural network comprising a 1D convolution on stacked intermediate logits obtained from auxiliary networks. The convolution operation reduces the integrates logits across layers, one class at a time. To this representation, we append uncertainty measures: confidence margin and entropy of each Auxiliary Network prediction. Using a linear projection on this intermediate representation, the *forecaster* provides a scalar weight $w_{(\psi^f)} \in [0,1]$ for each instance in the minibatch, which indicates the confidence of the *forecaster* in assessing whether the student network will correctly predict the input instance.

Post-hoc Adjustment of Forecaster Output. The forecaster $\mathcal{F}(z_f;\psi^f)$ outputs a weight $w_t \in [0,1]$. For a batch of input instances, forecaster output distribution can be skewed, affecting the stability of student model training and individual loss component contributions. Therefore, for a batch of samples (B), we employ the following post-hoc adjustment on the forecaster output logits (z_f) before utilizing the output for training the student, $w_{\mathrm{adj}} = \sigma\left(\varsigma^*\left(\frac{z_f - \mu_B}{\varsigma_B}\right) + \mu^*\right)$, where μ^* and ς^* are hyperparameters. Ultimately, the student is trained by minimizing, $\mathcal{L}_{\mathrm{tot}}^{(s)} = (1 - w_{\mathrm{adj}}) \cdot \mathcal{L}_{\mathrm{CE}}^{(s)} + w_{\mathrm{adj}} \cdot \mathcal{L}_{\mathrm{KD}}^{(s)}$

МЕТНОО	A	C	P	R	AVG.
KD (Hinton et al., 2015b)	52.7	49.7	72.1	74.2	62.2
KD_{+F}			70.5		
KD _{+F+ADJ} .	54.0	51.3	71.6	75.1	63.0

Table 1: **OOD** classification accuracy (%) on the **OfficeHome dataset**. KD denotes vanilla knowledge distillation, F incorporates the forecaster, and ADJ. further applies the proposed post-hoc adjustment to forecaster outputs. The adjustment ($\mu^*=0.5$ and $\varsigma^*=0.1$) consistently improves average performance across held-out domains.

4.3 ALGORITHM

We outline the distillation process with an interleaved training process where the student model, auxiliary networks and the forecaster are trained in an alternating pattern. For pre-determined k steps, the primary student model and the auxiliary networks are trained on the training domains using the training split of the training domains. Here, the forecaster remains frozen, and is used to generate instance-level weights to guide the student training. Subsequently, for next k steps, the forecaster is trained on the on a held-out validation split from the training data of the training domain(s) with the lastest checkpoints of the student model and the auxiliary networks, which are frozen. This process continues until the exit criterion of the algorithm is satisfied. Refer Algorithm 1 for an illustrated outline of the algorithm.

324 Algorithm 1 Interleaved Training of Student and Forecaster 325 326 **Require:** Student $\mathcal{M}_{\theta(S)}$, Teacher $\mathcal{M}_{\theta(T)}$, Forecaster $\mathcal{F}(\psi_f)$; $\mathcal{B}_{\text{train}}, \mathcal{B}_{\text{val}}$; T_s (student steps), T_f (forecaster steps), N (total steps), L (number of layers), μ^* (Adjustment mean), ς^* (Adjustment 327 328 1: Initialize $\boldsymbol{\theta}_0^{(S)}, \boldsymbol{\theta}_0^{(T)}, \boldsymbol{\psi}_0^f, [\varphi_0^{(\ell)}]_{\ell=1}^{L-1}$ 2: **for** t = 1 to N **do** 330 if $t \mod (T_s + T_f) < T_s$ then 331 ► Student and Auxiliary Network Update 4: 332 $\begin{array}{l} \text{Set } \mathcal{M}_{\boldsymbol{\theta}^{(S)}}, \mathcal{M}_{\boldsymbol{\theta}^{(T)}}, \mathcal{F}(\boldsymbol{\psi}_f), [\mathbf{A}_{\text{aux}}^{(\ell)}]_{\ell=1}^{L-1} \text{ to eval} \\ \mathbf{X} \leftarrow \text{concat}(\mathbf{x} \mid (\mathbf{x}, \mathbf{y}) \in \mathcal{B}_{train}) \end{array}$ 5: 333 6: 334 $[\mathbf{A}_{\mathrm{aux}}^{(\ell)}(\mathbf{X};\,\boldsymbol{\varphi})]_{\ell=1}^{L-1}, \hat{\mathbf{Y}} \leftarrow \mathcal{M}_{\boldsymbol{\theta}^{(S)}}(\mathbf{X})$ 7: 335 $\begin{aligned} &[\mathbf{A}_{\mathsf{aux}}(\mathbf{A}, \boldsymbol{\varphi})]_{\ell=1}, \mathbf{I} \leftarrow \mathcal{M}_{\boldsymbol{\theta}^{(S)}}(\mathbf{A}) \\ &\mathbf{O} \leftarrow \mathcal{F}([\mathbf{A}_{\mathsf{aux}}^{(\ell)}(\mathbf{X}; \boldsymbol{\varphi})]_{\ell=1}^{L-1}) \\ &\mu_B = \mathbf{O}.\mathsf{mean}\left(\right), \varsigma_B = \mathbf{O}.\mathsf{std}\left(\right) \\ &\mathbf{W} = \sigma\Big(\varsigma^*\left(\frac{\mathbf{O}-\mu_B}{\varsigma_B}\right) + \mu^*\Big) \\ &\mathsf{Set}\,\mathcal{M}_{\boldsymbol{\theta}^{(S)}}\,\mathsf{to}\,\,\mathbf{train}; \\ &\boldsymbol{\theta}_t^{(S)} \leftarrow \boldsymbol{\theta}_{t-1}^{(S)} - \frac{\eta}{|\mathcal{B}_{\mathsf{train}}|} \sum_{\mathbf{x} \in \mathcal{B}_{\mathsf{train}}} \mathcal{L}_{\mathsf{tot}}^{(s)}(\mathbf{x}; \boldsymbol{\theta}_{t-1}^{(S)}, \mathbf{W}) \\ &\mathbf{for}\,\,\ell = 1\,\mathsf{to}\,\,L - 1\,\mathbf{do} \end{aligned}$ 336 8: 337 9: 338 10: 339 11: 340 12: 341 13: 342 Set $\mathbf{A}_{\mathsf{aux}}^{(\ell)}$ to **train**; $\varphi_t^{(\ell)} = \varphi_{t-1}^{(\ell)} - \frac{\eta}{|\mathcal{B}_{\mathsf{train}}|} \sum_{\mathbf{x} \in \mathcal{B}_{\mathsf{train}}} \mathcal{L}_{\mathsf{aux}}^{(\ell)}(\mathbf{A}_{\mathsf{aux}}^{(\ell)}(\mathbf{x}); \varphi_{t-1}^{(\ell)})$ 14: 343 15: 344 16: 345 else 17: ► Forecaster Update 18: 347 Set $\mathcal{F}(\psi_f)$ to **train**, $\mathcal{M}_{\boldsymbol{\theta}^{(S)}}, [\mathbf{A}_{\mathsf{aux}}^{(\ell)}(\bullet; \varphi)]_{\ell=1}^{L-1}$ to **eval** $\psi_t^f \leftarrow \psi_{t-1}^f - \frac{\eta_f}{|\mathcal{B}_{\mathsf{val}}|} \sum_{\mathcal{B}_{\mathsf{val}}} \mathcal{L}^{(f)}(\theta^{(S)}, \psi_{t-1}^f)$ 19: 348 349 20: end if 350 21: **22: end for** 351

5 EXPERIMENTS

352 353

354 355

356

357

358

359

360 361

362

364

365

366

367

368

369

370

371

372

373

374

375

376

377

Evaluation. For the image-classification experiments, we leverage the DomainBed Suite (Gulrajani & Lopez-Paz, 2020). DomainBed provides a fair, standardized and reproducible setup for evaluating and comparing our method against best performing subset from a wide range of DG algorithms including ERM (Vapnik, 1998), GroupDRO (Sagawa et al., 2020), Mixup (Yan et al., 2020), MLDG (Li et al., 2017b), CORAL (Sun & Saenko, 2016), MMD (Li et al., 2018a), DANN (Ganin et al., 2016) and C-DANN (Li et al., 2018b). We use classification accuracy as the primary metric.

Model Pool. In our KD experiments, we distill from a robust teacher model. We experiment with both ResNet (He et al., 2015) and Vision Transformer (Dosovitskiy et al., 2021) as the teacher network. Specifically, we first identify the best-performing domain generalization algorithm at the teacher level per dataset using ResNet-152 and ViT-L/16, and use the best teacher network in the distillation process. For the student network, we select the smallest capacity model from the ResNet family: ResNet-18.

DATASET	DOMAINS	INSTANCES	LABELS
OFFICEHOME	4	15,588	65
PACS	4	9,991	7
VLCS	4	10,729	5
TERRAINCOGNITA	. 4	24,778	10

Table 2: Dataset statistics.

Datasets. We evaluate and compare our method on four benchmark datasets, including OfficeHome (Venkateswara et al., 2017), PACS (Li et al., 2017a), VLCS (Fang et al., 2013) and TerraIncognita (Beery et al., 2018). Table 2 provides the details of the domains, instances, and label counts for each dataset.

Baselines. We compare our method against three baselines. We include the Empirical Risk Minimization (ERM) (Vapnik, 1998) as a canonical domain generalization algorithm, supported by findings in Gulrajani & Lopez-Paz (2020). Then, we select the best-performing OOD algorithm iden-

tified at the teacher level as the second baseline. In this case, we retrain the student model without distillation. Finally, we include vanilla KD (Hinton et al., 2015b), where the student us trained to mimic the logits of the teacher without any further modifications.

Training Protocol. For a dataset comprising N domains, we adopt the leave-one-domain-out setup where the network is trained on N-1 domains and evaluated on the held-out domain. We reserve 20% from each of training domain as a validation set for model selection. In our method, remaining 80% training split is further divided, with 90% used to train the student backbone and auxiliary networks, and 10% held out to train the forecaster.

Auxiliary Network Design. For the ResNet student, each residual block (Total: 4) stage yields feature maps that encode progressively abstract representations of the input (He et al., 2015). To obtain class-wise representations from these intermediate features, we employ lightweight auxiliary heads following the design principles of prior works on multi-layer feature supervision (Szegedy et al., 2015; Lee et al., 2015) and early-exits (Xin et al., 2020; Liu et al., 2020). Each auxiliary head comprises two components: (i) a pooling operation to reduce the spatial dimensionality, and (ii) a feed-forward projection layer that maps the pooled features to a vector of dimension equal to the number of target classes.

Метнор	ARCH.	ОfficeHome	PACS	VLCS	TERRA.	AVG.
ERM (Vapnik, 1998) DB. SOTA	ResNet-18 ResNet-18	58.4 60.2	79.6 80.1	71.6 71.9	43.4 42.6	63.3 63.7
TEACHER NETWORK: Re	sNet-152					
KD (Hinton et al., 2015b)	ResNet-18	62.2	82.4	73.5	43.6	65.4
$KD_{+F+ADJ.}(Ours)$	ResNet-18	63.0	82.8	74.7	45.1	66.4
TEACHER NETWORK: Vi	T-L/16					
KD (Hinton et al., 2015b)	ResNet-18	63.7	81.4	75.2	40.6	65.2
$KD_{+F+ADJ.}$ (Ours)	ResNet-18	63.7	81.0	76.1	44.7	66.4

Table 3: Domain generalization accuracy (%) on four benchmark datasets. Best results are highlighted in bold and green. Here, DB. SOTA refers to the best-performing DG algorithm for the dataset, selected according to the evaluation protocol described in Appendix B.1.

6 RESULTS

Table 3 reports the domain generalization performance across four benchmark datasets. Here, ERM Vapnik (1998) corresponds to standard empirical risk minimization, where the network is trained solely on ground-truth supervision. Similarly, DB. SOTA refers to the strongest domain generalization algorithm among GroupDRO (Sagawa et al., 2020), Mixup (Yan et al., 2020), MLDG (Li et al., 2017b), CORAL (Sun & Saenko, 2016), MMD (Li et al., 2018a), DANN (Ganin et al., 2016) and C-DANN (Li et al., 2018b), selected as per findings described in Appendix B.1. In the KD (Hinton et al., 2015a) setting, the student is trained with an additional supervision from a larger network fine-tuned using the best-performing DG algorithm.

We observe that KD consistently outperforms both, the canonical empirical risk minimization (ERM) and the dataset-specific best-performing DG algorithm. This highlights the effectiveness of knowledge transfer from a robust teacher network. Building on the strong KD baseline, we evaluate our adaptive distillation setup, which augments KD with the forecaster meta-network that utilizes early readout signals from the student network. Our method improves over KD under both, a convolutional ResNet teacher and transformer based ViT teacher. With a robust ResNet-152 teacher, our approach achieves an average OOD accuracy of 66.4%, outperforming KD by +1.0%. Similarly, with a robust ViT-L/16 teacher, our method surpasses KD by +1.2%. These consistent improvements across datasets showcase our approach as a principled extention to the standard KD problem for domain generalization.

The Need For Post-hoc Adjustment. As shown in Table 1, in the absence output adjustment, adaptive KD with forecaster yields lower OOD performance as compared to vanilla KD. This motivates us to understand why an unadjusted forecaster may fail. We therefore analyze the dynamics

of the forecaster outputs during the distillation process with and without post-hoc adjustment to the forecaster outputs. Figure 3 showcases the change in forecaster outputs for a minibatch as training progresses. Without any correction to the forecaster output, the distribution gradually drifts towards extreme values close to 1. This drift can be attributed to the nature of the training of the forecaster. The forecaster is trained as a binary classifier to predict correctness of student network based on early readout predictions and uncertainty signals. Naturally, as the student improves, an increasing fraction of samples become easy, pushing forecaster towards overconfident predictions. This causes the weight to concentrate more to \mathcal{L}_{KD} as the training progresses, making model overfit to teacher's supervisory signals, affecting its generalizability.

The Role of Uncertainty Signals. We also analyze the impact of early layer uncertainty features on forecaster quality. As shown in Figure 2, including entropy $\mathcal{H}(p)$ and confidence margin (ε) from auxiliary predictions substantially improves the forecaster's predictive ability to determine easy and hard samples, with an increase in AUC by +26%. This demonstrates that early readouts provide rich signals of sample ambiguity, which can be leveraged by the forecaster, to effectively modulate the loss weighing. Together, these emipirical ablations on the forecaster confirm that both post-hoc adjustment and statistical uncertainty-aware design choice is essential to the forecaster meta-network.

7 CONCLUSION

In this work, we introduced a *forecaster* based re-weighing approach to standard offline knowledge distillation setting, where the *forecaster* leverages early layer readouts from the student model to adaptively modulate the distillation objective, resulting in an improved generalizability of the student network as opposed to vanilla KD. Our experiments across multiple benchmarks demonstrate the efficacy of our approach in improving OOD generalization of student network over vanilla KD baseline and canonical DG algorithms. Further, we provide insights on critical design choices for the *forecaster*: (1) post-hoc adjustment, which prevents *forecaster* collapse that results in overconfident predictions, and (2) Use of uncertainty measures such as entropy and confidence margin, which significantly improve *forecaster's* ability to distinguish between easy and difficult samples. Together, we establish our framework as a novel extension to standard offline KD, allowing robust generalization to unseen domains, from a student-centric design choice.

REFERENCES

- Guillaume Alain and Yoshua Bengio. Understanding intermediate layers using linear classifier probes. *arXiv preprint arXiv:1610.01644*, 2016.
- Görkem Algan and Ilkay Ulusoy. Meta soft label generation for noisy labels, 2021. URL https://arxiv.org/abs/2007.05836.
- Robert Baldock, Hartmut Maennel, and Behnam Neyshabur. Deep learning through the lens of example difficulty. *Advances in Neural Information Processing Systems*, 34:10876–10889, 2021.
- Sara Beery, Grant van Horn, and Pietro Perona. Recognition in terra incognita, 2018. URL https://arxiv.org/abs/1807.04975.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale, 2021. URL https://arxiv.org/abs/2010.11929.
- Chen Fang, Ye Xu, and Daniel N. Rockmore. Unbiased metric learning: On the utilization of multiple datasets and web images for softening bias. In *2013 IEEE International Conference on Computer Vision*, pp. 1657–1664, 2013. doi: 10.1109/ICCV.2013.208.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks, 2017. URL https://arxiv.org/abs/1703.03400.
- Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks, 2016. URL https://arxiv.org/abs/1505.07818.

- Ishaan Gulrajani and David Lopez-Paz. In search of lost domain generalization, 2020. URL https://arxiv.org/abs/2007.01434.
 - Alperen Görmez, Venkat R. Dasari, and Erdem Koyuncu. E2cm: Early exit via class means for efficient supervised and unsupervised learning. In 2022 International Joint Conference on Neural Networks (IJCNN), pp. 1–8, 2022. doi: 10.1109/IJCNN55064.2022.9891952.
 - Yizeng Han, Gao Huang, Shiji Song, Le Yang, Honghui Wang, and Yulin Wang. Dynamic neural networks: A survey, 2021. URL https://arxiv.org/abs/2102.04906.
 - Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015. URL https://arxiv.org/abs/1512.03385.
 - Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv* preprint arXiv:1503.02531, 2015a.
 - Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network, 2015b. URL https://arxiv.org/abs/1503.02531.
 - Rui Huang, Andrew Geng, and Yixuan Li. On the importance of gradients for detecting distributional shifts in the wild. *Advances in Neural Information Processing Systems*, 34:677–689, 2021.
 - Zeyi Huang, Andy Zhou, Zijian Lin, Mu Cai, Haohan Wang, and Yong Jae Lee. A sentence speaks a thousand images: Domain generalization through distilling clip with language guidance, 2023. URL https://arxiv.org/abs/2309.12530.
 - Fotis Iliopoulos, Vasilis Kontonis, Cenk Baykal, Gaurav Menghani, Khoa Trinh, and Erik Vee. Weighted distillation with unlabeled examples, 2022. URL https://arxiv.org/abs/2210.06711.
 - Nishant Jain, Karthikeyan Shanmugam, and Pradeep Shenoy. Selective classification using a robust meta-learning approach, 2024. URL https://arxiv.org/abs/2212.05987.
 - Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanas Phillips, Irena Gao, et al. Wilds: A benchmark of in-the-wild distribution shifts. In *International conference on machine learning*, pp. 5637–5664. PMLR, 2021.
 - Stefanos Laskaridis, Alexandros Kouris, and Nicholas D. Lane. Adaptive inference through early-exit networks: Design, challenges and directions. In *Proceedings of the 5th International Workshop on Embedded and Mobile Deep Learning*, MobiSys '21, pp. 1–6. ACM, June 2021. doi: 10.1145/3469116.3470012. URL http://dx.doi.org/10.1145/3469116.3470012.
 - Chen-Yu Lee, Saining Xie, Patrick Gallagher, Zhengyou Zhang, and Zhuowen Tu. Deeply-supervised nets. In *Artificial intelligence and statistics*, pp. 562–570. Pmlr, 2015.
 - Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M. Hospedales. Deeper, broader and artier domain generalization, 2017a. URL https://arxiv.org/abs/1710.03077.
 - Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M. Hospedales. Learning to generalize: Meta-learning for domain generalization, 2017b. URL https://arxiv.org/abs/1710.03463.
 - Haoliang Li, Sinno Jialin Pan, Shiqi Wang, and Alex C. Kot. Domain generalization with adversarial feature learning. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5400–5409, 2018a. doi: 10.1109/CVPR.2018.00566.
 - Xin-Chun Li, Wen-Shu Fan, Shaoming Song, Yinchuan Li, Shao Yunfeng, De-Chuan Zhan, et al. Asymmetric temperature scaling makes larger networks teach well again. *Advances in neural information processing systems*, 35:3830–3842, 2022.
 - Ya Li, Mingming Gong, Xinmei Tian, Tongliang Liu, and Dacheng Tao. Domain generalization via conditional invariant representation, 2018b. URL https://arxiv.org/abs/1807.08479.

- Weijie Liu, Peng Zhou, Zhiruo Wang, Zhe Zhao, Haotang Deng, and Qi Ju. FastBERT: a self-distilling BERT with adaptive inference time. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (eds.), *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 6035–6044, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.537. URL https://aclanthology.org/2020.acl-main.537/.
 - Raphael Gontijo Lopes, Stefano Fenu, and Thad Starner. Data-free knowledge distillation for deep neural networks, 2017. URL https://arxiv.org/abs/1710.07535.
 - Yoshitomo Matsubara, Marco Levorato, and Francesco Restuccia. Split computing and early exiting for deep learning applications: Survey and research challenges. *ACM Computing Surveys*, 55 (5):1–30, December 2022. ISSN 1557-7341. doi: 10.1145/3527155. URL http://dx.doi.org/10.1145/3527155.
 - Omar Montasser, Han Shao, and Emmanuel Abbe. Transformation-invariant learning and theoretical guarantees for ood generalization. *Advances in Neural Information Processing Systems*, 37: 108649–108673, 2024.
 - Anshul Nasery, Sravanti Addepalli, Praneeth Netrapalli, and Prateek Jain. Daft: Distilling adversarially fine-tuned models for better ood generalization, 2022. URL https://arxiv.org/abs/2208.09139.
 - Aniruddh Raghu, Maithra Raghu, Simon Kornblith, David Duvenaud, and Geoffrey Hinton. Teaching with commentaries, 2021. URL https://arxiv.org/abs/2011.03037.
 - Mengye Ren, Wenyuan Zeng, Bin Yang, and Raquel Urtasun. Learning to reweight examples for robust deep learning, 2019. URL https://arxiv.org/abs/1803.09050.
 - Alexander Robey, George J Pappas, and Hamed Hassani. Model-based domain generalization. *Advances in Neural Information Processing Systems*, 34:20210–20229, 2021.
 - Shiori Sagawa, Pang Wei Koh, Tatsunori B. Hashimoto, and Percy Liang. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization, 2020. URL https://arxiv.org/abs/1911.08731.
 - Shreyas Saxena, Oncel Tuzel, and Dennis DeCoste. Data parameters: A new family of parameters for learning a differentiable curriculum. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/file/926ffc0ca56636b9e73c565cf994ea5a-Paper.pdf.
 - Roy Schwartz, Gabriel Stanovsky, Swabha Swayamdipta, Jesse Dodge, and Noah A. Smith. The right tool for the job: Matching model and instance complexities. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (eds.), *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 6640–6651, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.593. URL https://aclanthology.org/2020.acl-main.593/.
 - Jun Shu, Qi Xie, Lixuan Yi, Qian Zhao, Sanping Zhou, Zongben Xu, and Deyu Meng. Meta-weightnet: Learning an explicit mapping for sample weighting, 2019. URL https://arxiv.org/abs/1902.07379.
 - Durga Sivasubramanian, Ayush Maheshwari, Pradeep Shenoy, Prathosh AP, and Ganesh Ramakrishnan. Adaptive mixing of auxiliary losses in supervised learning, 2022. URL https://arxiv.org/abs/2202.03250.
 - Rupesh K Srivastava, Klaus Greff, and Jürgen Schmidhuber. Training very deep networks. *Advances in neural information processing systems*, 28, 2015.
 - Baochen Sun and Kate Saenko. Deep coral: Correlation alignment for deep domain adaptation, 2016. URL https://arxiv.org/abs/1607.01719.

- Tianxiang Sun, Xiangyang Liu, Wei Zhu, Zhichao Geng, Lingling Wu, Yilong He, Yuan Ni, Guotong Xie, Xuanjing Huang, and Xipeng Qiu. A simple hash-based early exiting approach for language understanding and generation, 2022. URL https://arxiv.org/abs/2203.01670.
 - Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1–9, 2015.
 - Rishabh Tiwari, Durga Sivasubramanian, Anmol Mekala, Ganesh Ramakrishnan, and Pradeep Shenoy. Using early readouts to mediate featural bias in distillation, 2023. URL https://arxiv.org/abs/2310.18590.
 - Rishabh Tiwari, Durga Sivasubramanian, Anmol Mekala, Ganesh Ramakrishnan, and Pradeep Shenoy. Using early readouts to mediate featural bias in distillation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 2638–2647, 2024.
 - V.N. Vapnik. *Statistical Learning Theory*. Adaptive and learning systems for signal processing, communications, and control. Wiley, 1998. ISBN 9788126528929. URL https://books.google.co.in/books?id=RWrlkQEACAAJ.
 - Hemanth Venkateswara, Jose Eusebio, Shayok Chakraborty, and Sethuraman Panchanathan. Deep hashing network for unsupervised domain adaptation, 2017. URL https://arxiv.org/abs/1706.07522.
 - Yoav Wald, Amir Feder, Daniel Greenfeld, and Uri Shalit. On calibration and out-of-domain generalization. Advances in neural information processing systems, 34:2215–2227, 2021.
 - Yisen Wang, Weiyang Liu, Xingjun Ma, James Bailey, Hongyuan Zha, Le Song, and Shu-Tao Xia. Iterative learning with open-set noisy labels, 2018. URL https://arxiv.org/abs/1804.00092.
 - Yufei Wang, Haoliang Li, Lap pui Chau, and Alex C. Kot. Embracing the dark knowledge: Domain generalization using regularized knowledge distillation, 2021. URL https://arxiv.org/abs/2107.02629.
 - Ji Xin, Raphael Tang, Jaejun Lee, Yaoliang Yu, and Jimmy Lin. DeeBERT: Dynamic early exiting for accelerating BERT inference. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (eds.), *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 2246–2251, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.204. URL https://aclanthology.org/2020.acl-main.204/.
 - Canwen Xu and Julian McAuley. A survey on dynamic neural networks for natural language processing, 2023. URL https://arxiv.org/abs/2202.07101.
 - Shen Yan, Huan Song, Nanxiang Li, Lincan Zou, and Liu Ren. Improve unsupervised domain adaptation with mixup training, 2020. URL https://arxiv.org/abs/2001.00677.
 - Xinli Yue, Mou Ningping, Qian Wang, and Lingchen Zhao. Revisiting adversarial robustness distillation from the perspective of robust fairness. *Advances in Neural Information Processing Systems*, 36:30390–30401, 2023.
 - Di Zhao, Jingfeng Zhang, Hongsheng Hu, Philippe Fournier-Viger, Gillian Dobbie, and Yun Sing Koh. Balancing invariant and specific knowledge for domain generalization with online knowledge distillation. In James Kwok (ed.), *Proceedings of the Thirty-Fourth International Joint Conference on Artificial Intelligence, IJCAI-25*, pp. 2440–2448. International Joint Conferences on Artificial Intelligence Organization, 8 2025. doi: 10.24963/ijcai.2025/272. URL https://doi.org/10.24963/ijcai.2025/272. Main Track.
 - Haoran Zhao, Xin Sun, Junyu Dong, Zihe Dong, and Qiong Li. Knowledge distillation via instance-level sequence learning, 2021. URL https://arxiv.org/abs/2106.10885.

 Shanshan Zhao, Mingming Gong, Tongliang Liu, Huan Fu, and Dacheng Tao. Domain generalization via entropy regularization. *Advances in neural information processing systems*, 33:16096–16107, 2020.

- Andy Zhou, Jindong Wang, Yu-Xiong Wang, and Haohan Wang. Distilling out-of-distribution robustness from vision-language foundation models. *Advances in Neural Information Processing Systems*, 36:32938–32957, 2023.
- Kaiyang Zhou, Yuanhan Zhang, Yuhang Zang, Jingkang Yang, Chen Change Loy, and Ziwei Liu. On-device domain generalization. *arXiv preprint arXiv:2209.07521*, 2022.
- Wangchunshu Zhou, Canwen Xu, Tao Ge, Julian McAuley, Ke Xu, and Furu Wei. Bert loses patience: Fast and robust inference with early exit, 2020. URL https://arxiv.org/abs/2006.04152.

A IMPLEMENTATION AND HYPERPARAMETERS

For a fair comparison and reproducibility, we follow the hyperparameter settings reported in Gulrajani & Lopez-Paz (2020) and use their released codebase as the foundation of our implementation.

B Additional Experimental Results

B.1 SELECTION OF ROBUST DG ALGORITHM

ALGORITHM	Оғысеноме	VLCS	PACS	TERRAINC.	AVG.
ERM (Vapnik, 1998)	68.4	77.6	86.4	48.6	70.3
GROUPDRO (Sagawa et al., 2020)	68.6	78.0	87.0	46.7	70.1
MIXUP (Yan et al., 2020)	70.4	78.0	87.5	46.1	70.5
MLDG (Li et al., 2017b)	56.3	72.1	68.1	33.5	57.5
CORAL (Sun & Saenko, 2016)	68.6	77.1	87.7	48.4	70.5
MMD (Li et al., 2018a)	69.2	77.5	85.6	48.7	70.3
DANN (Ganin et al., 2016)	69.7	79.7	86.3	47.4	70.8
CDANN (Li et al., 2018b)	70.0	79.0	86.0	48.3	70.8

Table 4: Average OOD classification accuracies (%) for all datasets with ResNet-152. Model selection: training-domain validation set.

ALGORITHM	OFFICEHOME	VLCS	PACS	TERRAINC.	Avg.
ERM (Vapnik, 1998)	71.6	78.4	83.9	38.8	68.2
GROUPDRO (Sagawa et al., 2020)	71.5	77.3	80.3	33.7	65.7
MIXUP (Yan et al., 2020)	72.2	76.7	81.6	42.5	68.3
MLDG (Li et al., 2017b)	70.5	76.3	81.3	40.0	67.0
CORAL (Sun & Saenko, 2016)	71.9	78.2	84.4	40.8	68.8
MMD (Li et al., 2018a)	71.3	77.5	81.9	38.8	67.4
DANN (Ganin et al., 2016)	72.1	78.7	84.3	43.1	69.6
CDANN (Li et al., 2018b)	71.4	77.2	78.6	42.9	67.5

Table 5: Average OOD classification accuracies (%) for all datasets with ViT-L/16. Model selection: training-domain validation set.

B.2 Domain-wise Results on Benchmark Datasets

Tables 6-13 report quantitative figures for domain generalization on the OfficeHome (Venkateswara et al., 2017), PACS (Li et al., 2017a), VLCS (Fang et al., 2013), and TerraIncognita (Beery et al., 2018) datasets. We use the ResNet-18 as the primary network in all experiments. Tables 6-9 use ResNet-152 as the teacher network, while Tables 10-13 use ViT-L/16 as the teacher network. Each column in the table represents the held-out domain.

Метнор	A	С	P	R	AVG.
ERM (Vapnik, 1998)	50.4	48.3	65.4	69.4	58.4
MIXUP (Yan et al., 2020)	51.5	48.6	70.3	70.2	60.2
KD (Hinton et al., 2015b)	52.7	49.7	72.1	74.2	62.2
$KD_{+F+ADJ.}$	54.0	51.3	71.6	75.1	63.0

Table 6: OOD classification accuracy (%) on the OfficeHome dataset. Model selection: training-domain validation set. KD experiments use a ResNet-152 teacher network.

Метнор	A	С	P	S	AVG.
ERM (Vapnik, 1998)	77.2	73.3	94.7	73.2	79.6
CORAL (Sun & Saenko, 2016)	77.3	73.9	94.7	74.5	80.1
KD (Hinton et al., 2015b)	82.4	76.4	95.0	75.9	82.4
$K_{D+F+ADJ}$.	82.4	77.0	95.7	76.0	82.8

Table 7: OOD classification accuracy (%) on the PACS dataset. Model selection: training-domain validation set. KD experiments use a ResNet-152 teacher network.

Метнор	С	S	L	V	AVG.
ERM (Vapnik, 1998)	96.6	58.5	65.4	65.9	71.6
DANN (Ganin et al., 2016)	96.9	58.8	65.7	66.2	71.9
KD (Hinton et al., 2015b)	97.2	60.0	67.9	68.8	73.5
$KD_{+F+ADJ.}$	98.3	61.3	69.8	69.2	74.7

Table 8: OOD classification accuracy (%) on the VLCS dataset. Model selection: training-domain validation set. KD experiments use a ResNet-152 teacher network.

Метнор	L100	L38	L43	L46	AVG.
ERM (Vapnik, 1998)	53.2	33.0	51.4	35.9	43.4
MMD (Li et al., 2018a)	47.4	37.4	50.0	35.5	42.6
KD (Hinton et al., 2015b)	50.3	36.8	52.7	34.6	43.6
$KD_{+F+ADJ.}$	56.7	33.6	54.4	35.5	45.1

Table 9: OOD classification accuracy (%) on the TerraIncognita dataset. Model selection: training-domain validation set. KD experiments use a ResNet-152 teacher network.

МЕТНОО	A	С	P	R	AVG.
ERM (Vapnik, 1998)	50.4	48.3	65.4	69.4	58.4
MIXUP (Yan et al., 2020)	51.5	48.6	70.3	70.2	60.2
KD (Hinton et al., 2015b)	55.4	53.0	72.6	74.0	63.7
$KD_{+F+ADJ.}$	56.3	52.6	72.5	73.2	63.7

Table 10: OOD classification accuracy (%) on the OfficeHome dataset. Model selection: training-domain validation set. KD experiments use a ViT-L/16 teacher network.

Метнор	A	С	P	S	AVG.
ERM (Vapnik, 1998)	77.2	73.3	94.7	73.2	79.6
CORAL (Sun & Saenko, 2016)	77.3	73.9	94.7	74.5	80.1
KD (Hinton et al., 2015b)	79.7	75.6	95.9	74.6	81.4
$K_{D+F+ADJ.}$	78.3	74.7	95.9	75.0	81.0

Table 11: OOD classification accuracy (%) on the PACS dataset. Model selection: training-domain validation set. KD experiments use a ViT-L/16 teacher network.

Метнор	С	S	L	V	AVG.
ERM (Vapnik, 1998)	96.6	58.5	65.4	65.9	71.6
DANN (Ganin et al., 2016)	96.9	58.8	65.7	66.2	71.9
KD (Hinton et al., 2015b)	95.9	61.5	69.7	73.6	75.2
$K_{D_{+F+ADJ.}}$	96.1	62.2	71.2	74.8	76.1

Table 12: OOD classification accuracy (%) on the VLCS dataset. Model selection: training-domain validation set. KD experiments use a ViT-L/16 teacher network.

МЕТНОО	L100	L38	L43	L46	AVG.
ERM (Vapnik, 1998)	53.2	33.0	51.4	35.9	43.4
MMD (Li et al., 2018a)	47.4	37.4	50.0	35.5	42.6
KD (Hinton et al., 2015b)	47.4	32.1	50.2	32.7	40.6
OURS	55.7	36.8	52.4	34.2	44.7

Table 13: OOD classification accuracy (%) on the TerraIncognita dataset. Model selection: training-domain validation set. KD experiments use a ViT-L/16 teacher network.