
Robust Calibration with Multi-domain Temperature Scaling

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Uncertainty quantification is essential for the reliable deployment of machine
2 learning models to high-stakes application domains. Uncertainty quantification is
3 all the more challenging when training distribution and test distribution are different,
4 even the distribution shifts are mild. Despite the ubiquity of distribution shifts
5 in real-world applications, existing uncertainty quantification approaches mainly
6 study the in-distribution setting where the train and test distributions are the same.
7 In this paper, we develop a systematic calibration model to handle distribution shifts
8 by leveraging data from multiple domains. Our proposed method—multi-domain
9 temperature scaling—uses the heterogeneity in the domains to improve calibration
10 robustness under distribution shift. Through experiments on three benchmark data
11 sets, we find our proposed method outperforms existing methods as measured on
12 both in-distribution and out-of-distribution test sets.

13 1 Introduction

14 To make learning systems reliable and fault-tolerant, predictions must be accompanied by uncertainty
15 estimates. A significant challenge to accurately codifying uncertainty is the distribution shift that
16 typically arises over the course of a system’s deployment [31]. For example, suppose health providers
17 from 20 different hospitals employ a model to make diagnostic predictions from fMRI data. The
18 distributions across hospitals could be quite different as a result of differing patient populations,
19 machine conditions, and so on. In such a setting, it is critical to provide uncertainty quantification that
20 is valid for *every* hospital—not just on average across all hospitals. Going even further, our uncertainty
21 quantification should be informative when a new 21st hospital goes online, even if the distribution
22 shifts from those already encountered. In this work, we study calibration in the multi-domain setting.
23 We find that by requiring accurate calibration across all observed domains, our method provides more
24 accurate uncertainty quantification on unseen domains.

25 Calibration is a core topic in learning [29, 24, 8, 20, 11, 2], but most techniques are targeted at settings
26 with no distribution shift. To see this, we consider a simple experiment on the ImageNet-C [13]
27 dataset, which consists of 76 domains. Here, each domain corresponds to one type of data corruption
28 applied with a certain severity. We apply the temperature scaling technique [11] on the pooled data
29 from all domains. In Figure 1(a) and 1(b), we display the reliability diagrams for the pooled data
30 and for one individual domain. We find that even under a relatively mild distribution shift—i.e.,
31 subpopulation shift from the mixture of all domains to the single domain—temperature scaling does
32 not produce calibrated confidence estimates on the stand-alone domain. This behavior is pervasive;
33 in Figure 1(c), we see that the calibration on individual domains is much worse than the the reliability
34 diagram from the pooled data would suggest.

35 To address this issue, we develop a new algorithm, multi-domain temperature scaling, that leverages
36 multi-domain structure in the data. Our algorithm takes a base model and learns a calibration function

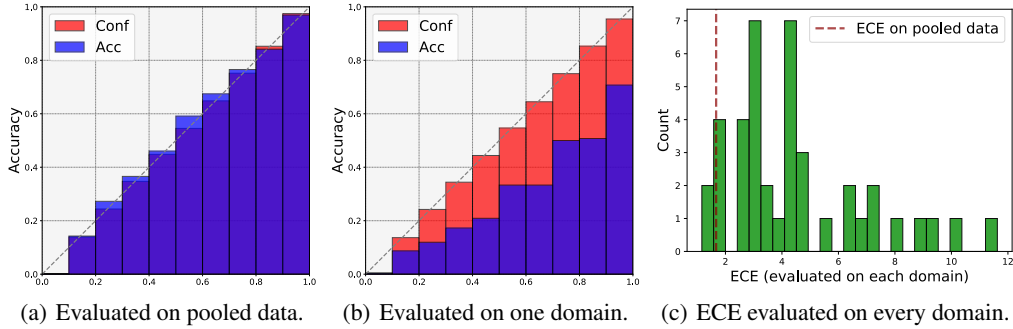


Figure 1: Reliability diagrams and expected calibration error histograms for temperature scaling with a ResNet-50 on ImageNet-C. We use temperature scaling to obtain adjusted confidences for the ResNet-50 model. **(a)** Reliability diagram evaluated on the pooled data of ImageNet-C. **(b)** Reliability diagram evaluated on data from one domain (Gaussian corruption with severity 5) in ImageNet-C. **(c)** Calibration evaluated on every domain in ImageNet-C as well as the pooled ImageNet-C (measured in ECE, lower is better).

37 that maps each input to a different temperature parameter that is used for adjusting confidence in the
 38 base model. Empirically, we find our algorithm significantly outperforms temperature scaling on
 39 three real-world multi-domain datasets. In particular, in contrast to temperature scaling, our proposed
 40 algorithm is able to provide well-calibrated confidence on each domain. Moreover, our algorithm
 41 largely improves robustness of calibration under distribution shifts. This is expected, because if the
 42 calibration method performs well on every domain, it is likely to have learned some structure that
 43 generalizes to unseen domains. Theoretically, we analyze the multi-domain calibration problem in
 44 the regression setting, providing guidance about the conditions under which robust calibration is
 45 possible.

46 **Contributions.** The main contributions of our work are as follows: Algorithmically, we develop a
 47 new calibration method that generalizes the widely used temperature scaling concept from single-
 48 domain to multi-domain. The proposed new method exploits multi-domain structure in the data
 49 distribution, which enables model calibration on every domain. We conduct detailed experiments on
 50 three real-world multi-domain datasets and demonstrate that our method significantly outperforms
 51 existing calibration methods on *both in-distribution domains and unseen out-of-distribution domains*.
 52 Theoretically, we study multi-domain calibration in the regression setting and develop a theoretical
 53 understanding of robust calibration in this setting.

54 Related Work

55 **Calibration methods.** There is a large literature on calibrating the well-trained machine learning
 56 models, including histogram binning [40], isotonic regression [41], conformal prediction [35], Platt
 57 scaling [29], and temperature scaling [11]. These calibration methods apply a validation set and post-
 58 process the model outputs. As shown in Guo et al. [11], temperature scaling, a simple method that
 59 uses a single (temperature) parameter for rescaling the logits, performs surprisingly well on calibrating
 60 confidences for deep neural networks. We focus on this approach in our work. More broadly, there has
 61 been much recent work develop methods to improve calibration for deep learning models, including
 62 augmentation-based training [33, 16], self-supervised learning [15], ensembling [20], and Bayesian
 63 neural networks [8, 9], as well as statistical guarantees for calibration with black-box models [1].

64 **Calibration under distribution shifts.** Ovia et al. [25] conduct an empirical study on model
 65 calibration under distribution shifts and find that models are much less calibrated under distribution
 66 shifts. Minderer et al. [22] revisit calibration of recent state-of-the-art image classification models
 67 under distribution shifts and study the relationship between calibration and accuracy. Wald et al.
 68 [36] study model calibration and out-of-distribution generalization. Other works consider providing
 69 uncertainty estimates under structured distribution shifts, such as covariate shift [34, 27], label
 70 shift [30], and f -divergence balls [5]. Another line of work studies calibration in the domain
 71 adaptation setting [37, 26], which require unlabeled samples from the target domain.

72 2 Problem setup

73 **Notation.** We denote the input space and the label set by $\mathcal{X} \subseteq \mathbb{R}^d$ and $\mathcal{Y} = \{1, \dots, J\}$. We let $[x]_i$
 74 denote the i -th element of vector x . We use $\mathcal{P}(X)$ to denote the marginal feature distribution on input

75 space \mathcal{X} , $\mathcal{P}(Y|X)$ to denote the conditional distribution, and $\mathcal{P}(X, Y)$ to denote the joint distribution.
 76 For the multiple domains scenario, we let $\mathcal{P}_k(X)$ and $\mathcal{P}_k(Y|X)$ denote the feature distribution and
 77 conditional distribution for the k -th domain. We let $f : \mathcal{X} \rightarrow \mathbb{R}^J$ denote the base model, e.g., a deep
 78 neural network, where J is the total number of classes. We assume f returns an (unnormalized) vector
 79 of logits. Throughout the paper, the base model is trained with training data and will not be modified.
 80 The class prediction of model f on input $x \in \mathcal{X}$ is denoted by $\hat{y} = \operatorname{argmax}_{j \in \{1, \dots, J\}} [f(x; \theta)]_j$. We
 81 use $\mathbf{1}\{\cdot\}$ to represent the indicator function. We use $h(\cdot; f, \beta) : \mathcal{X} \rightarrow [0, 1]$ to denote a *calibration*
 82 *map* (parameterized by β) that takes an input $x \in \mathcal{X}$ and returns a confidence score—this is a
 83 post-processing of the base model f . We let $\hat{\pi} = h(x; f, \beta) \in [0, 1]$ denote the confidence estimate
 84 for sample x when using model f . For instance, if we have 100 predictions $\{\hat{y}_1, \dots, \hat{y}_{100}\}$ with
 85 confidence $\hat{\pi}_1 = \dots = \hat{\pi}_{100} = 0.7$, then the accuracy of f is expected to be 70% on these 100
 86 samples (if the confidence estimate is well calibrated). Data from the domains $\mathcal{P}_1, \dots, \mathcal{P}_K$ are used
 87 for learning the calibration models, and we call the *in-distribution* (InD) domains. We use $\tilde{\mathcal{P}}$ to denote
 88 the unseen *out-of-distribution* (OOD) domain which is not used for calibrating the base model. Our
 89 goal is to learn a calibration map h that is well calibrated on the OOD domain $\tilde{\mathcal{P}}$. To do this, we will
 90 learn a calibration map that does well on all InD domains simultaneously.

91 To measure calibration, we first review the definition of approximate expected calibration error.

92 **Definition 2.1** (ECE). For a set of samples $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ with $(x_i, y_i) \stackrel{\text{i.i.d.}}{\sim} \mathcal{P}(X, Y)$, the
 93 (empirical) expected calibration error (ECE) with M bins evaluated on \mathcal{D} is defined as

$$\text{ECE}(\mathcal{D}, M) = \sum_{m=1}^M \frac{|B_m|}{n} |\text{Acc}(B_m) - \text{Conf}(B_m)|, \quad (1)$$

94 and B_m , $\text{acc}(B_m)$, $\text{conf}(B_m)$ are defined as

$$B_m = \{i \in [n] : \hat{\pi}_i \in ((m-1)/M, m/M)\},$$

$$\text{Acc}(B_m) = (1/|B_m|) \sum_{i \in B_m} \mathbf{1}\{\hat{y}_i = y_i\}, \quad \text{Conf}(B_m) = (1/|B_m|) \sum_{i \in B_m} \hat{\pi}_i,$$

95 where $\hat{\pi}_i$ and \hat{y}_i are the confidence and predicted label of sample x_i .

96 The empirical ECE defined in Eq. (1) approximates the expected calibration error (ECE) $\mathbb{E}[|p - \mathbb{P}(\hat{y} =$
 97 $y | \hat{\pi} = p)|]$ with bin size equal to M [24, 11]; see [21] for statistical results about about the empirical
 98 ECE as an estimator. The perfect calibrated map corresponds to the case when $\mathbb{P}(\hat{y} = y | \hat{\pi} = p) = p$
 99 holds for all $p \in [0, 1]$.

100 **Multi-domain calibration.** Although the standard ECE measurement in Eq. (1) provides informative
 101 evaluations for various calibration methods in the single-domain scenario, it does not provide fine-
 102 grained evaluations when the dataset consists of multiple domains, $\mathcal{P}_1, \dots, \mathcal{P}_K$. It is possible that the
 103 ECE evaluated on the pooled data $\mathcal{D}_K^{\text{pooled}} = \mathcal{D}_1 \cup \dots \cup \mathcal{D}_K$ is small while the ECE evaluated on one
 104 of the domains is large. For example, as shown in Figure 1(c), there may exist a domain, $k \in [K]$,
 105 such that the ECE evaluated on domain k is much higher than the ECE evaluated on the pooled
 106 dataset, i.e., $\text{ECE}(\mathcal{D}_k) \gg \text{ECE}(\mathcal{D}_K^{\text{pooled}})$. In the fMRI application mentioned in Section 1, producing
 107 well-calibrated confidence on data from every hospital is a more desirable property compared to only
 108 being calibrated on the pooled data from all hospitals. Therefore, it is natural to consider the ECE
 109 evaluated on every domain, which we refer to as “per-domain ECE.” Next, we introduce the notion of
 110 Multi-domain ECE to formalize per-domain calibration.

111 **Definition 2.2** (Multi-domain ECE). For a dataset $\mathcal{D}_K^{\text{pooled}} = \mathcal{D}_1 \cup \dots \cup \mathcal{D}_K$ consisting of samples
 112 from K domains, where $\mathcal{D}_k = \{(x_{i,k}, y_{i,k})\}_{i=1}^{n_k}$ and $(x_{i,k}, y_{i,k}) \stackrel{\text{i.i.d.}}{\sim} \mathcal{P}_k(X, Y)$, the (empirical)
 113 multi-domain expected calibration error (Multi-domain ECE) with M bins evaluated on $\mathcal{D}_K^{\text{pooled}}$ is
 114 defined as $\text{MDECE}(\mathcal{D}_K^{\text{pooled}}) = \frac{1}{K} \sum_{k=1}^K \text{ECE}(\mathcal{D}_k)$.

115 Compared with the standard ECE evaluated on the pooled dataset, multi-domain ECE provides
 116 information about per-domain model calibration. In the multi-domain setting, we aim to learn a
 117 calibration map \hat{h} that can produce calibrated confidence estimates on every InD domain. Intuitively,
 118 if the unseen OOD domain $\tilde{\mathcal{D}}$ is similar to one or multiple InD domains, \hat{h} can still provide reliable
 119 confidence estimates on the new domain. We formally study the connection between “well-calibrated
 120 on each InD domain” and “robust calibration on the OOD domain” in Section 5.

121 **Temperature scaling.** Next, we review a simple and effective calibration method, named temperature
 122 scaling (TS) [29, 11], that is widely used in single-domain model calibration. Temperature scaling
 123 applies a single parameter $T > 0$ and produces the confidence prediction for the base model f as

$$h^{\text{ts}}(x; f, T) = \max_{j \in \{1, \dots, J\}} [\text{Softmax}(f(x)/T)]_j,$$

124 where $[\text{Softmax}(z)]_j = \exp([z]_j) / \sum_{i=1}^J \exp([z]_i)$. The parameter T is the so-called *temperature*,
 125 with larger temperature yielding more diffuse probability estimates. To learn the temperature
 126 parameter T from dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$, Guo et al. [11] propose to find T by solving the
 127 following convex optimization problem,

$$\min_T \mathcal{L}_{\text{TS}}(T) := - \sum_{i=1}^n \sum_{j=1}^J \mathbf{1}\{y_i = j\} \cdot \log([\text{Softmax}(f(x_i)/T)]_j), \quad (2)$$

128 which optimizes the temperature parameter such that the negative log likelihood is minimized. We use
 129 TS-Alg to denote the temperature scaling learning algorithm; given inputs dataset \mathcal{D} and base model
 130 f , TS-Alg outputs the learned temperature parameter by solving Eq. (2), e.g., $\hat{T} = \text{TS-Alg}(\mathcal{D}, f)$.

131 3 Multi-domain temperature scaling

132 We propose our algorithm—multi-domain temperature scaling—that aims to improve the calibration
 133 on each domain. One key observation is that if we apply temperature scaling to each domain
 134 separately, then TS is able to produce calibrated confidence on every domain. Therefore, the question
 135 becomes how to “aggregate” these temperature scaling models and learn one calibration model,
 136 denoted by \hat{h} , that has similar performance to the k -th calibration model \hat{h}_k evaluated on domain k
 137 for every $k \in [K]$.

138 At a high level, we propose to learn a calibration model that maps samples from the input space \mathcal{X} to
 139 the temperature space \mathbb{R}_+ . To start with, we learn the temperature parameter \hat{T}_k for the base model
 140 on every domain k by applying temperature scaling on \mathcal{D}_k . Next, we apply the base deep model to
 141 compute feature embeddings of samples from different domains,¹ and label feature embeddings from
 142 the k -th domain with \hat{T}_k . In particular, we construct K new datasets, $\hat{\mathcal{D}}_1, \dots, \hat{\mathcal{D}}_K$, where each dataset
 143 contains feature embeddings and temperature labels from one domain, i.e., $\hat{\mathcal{D}}_k = \{(\Psi(x_{i,k}), \hat{T}_k)\}_{i=1}^{n_k}$.
 144 Finally, we apply linear regression on these labeled datasets. In detail, our algorithm is as follows:

- 145 1. **Learn temperature scaling model for each domain.** For every domain k , we learn
 146 temperature \hat{T}_k by applying temperature scaling on validation data $\mathcal{D}_k = \{(x_{i,k}, y_{i,k})\}_{i=1}^{n_k}$
 147 from k -th domain, i.e., $\hat{T}_k = \text{TS-Alg}(\mathcal{D}_k, f)$ and TS-Alg denotes the TS algorithm.
- 148 2. **Learn linear regression of temperatures.** Extract the feature embeddings of the base deep
 149 model f on each domain. Use $\Psi(x_{i,k}) \in \mathbb{R}^p$ to denote the feature embedding of the i -th
 150 sample from k -th domain. Then we learn $\hat{\theta}$ by solving the following optimization problem,

$$\hat{\theta} = \underset{\theta}{\text{argmin}} \sum_{k=1}^K \sum_{i=1}^{n_k} \left(\langle \Psi(x_{i,k}), \theta \rangle - \hat{T}_k \right)^2.$$

- 151 3. **Predict temperature on unseen test samples.** Given an unseen test sample \tilde{x} , we first
 152 compute the predicted temperature \tilde{T} using the learned linear model $\tilde{T} = \langle \Psi(x_{i,k}), \hat{\theta} \rangle$. Then
 153 we output the confidence estimate for sample \tilde{x} as

$$\tilde{\pi} = \max_j \left[\text{Softmax}(f(\tilde{x})/\tilde{T}) \right]_j.$$

154 We denote our proposed method by MD-TS (**M**ult-**D**omain **T**emperature **S**caling). A presentation of
 155 the algorithm in pseudocode can be found in Algorithm 1, Appendix A.

156 We pause to consider the basic concept in more detail. The goal of our proposed algorithm is to
 157 predict the best temperature for samples from different several domains. In an ideal setting where

¹We use the penultimate layer outputs of model f as the feature embeddings by default.

158 the learned linear model $\hat{\theta}$ results in good calibration on *every* InD domain, we can expect that $\hat{\theta}$
 159 will continue to yield good calibration on the OOD domain $\tilde{\mathcal{P}}$ when $\tilde{\mathcal{P}}$ is close to one or several InD
 160 domains. For example, $\tilde{\mathcal{P}}$ will work well if $\tilde{\mathcal{P}}$ is a mixture of the K domains, i.e., $\tilde{\mathcal{P}} = \sum_{k=1}^K \alpha_k \mathcal{P}_k$
 161 and $\alpha \in \Delta^{K-1}$. Regarding the algorithmic design, linear regression is one of the simplest models
 162 for solving the regression problem. It is computationally fast to learn such linear models as well as
 163 make predictions on new samples, making it attractive. We test alternative, more flexible, regression
 164 algorithms in Section 4 but do not observe significant gains over linear regression.

165 To illustrate how our proposed algorithm MD-TS performs differently from standard TS, we return to the
 166 ImageNet-C dataset. We compare the predicted temperature of our algorithm on new samples from domain
 167 k with the temperature that results from running TS on domain k alone. The results are summarized in
 168 Figure 2, where each circle corresponds to the mean predicted temperature on one InD domain. For each
 169 domain, we also visualize the standard deviation of the predicted temperatures for samples from that domain
 170 (the horizontal bar around each point). We find that our algorithm predicts the temperature quite well. Note that
 171 it does not have access to the domain index information of the fresh samples. By contrast, TS always uses the
 172 same temperature, regardless of the input point.
 173
 174
 175
 176
 177
 178
 179

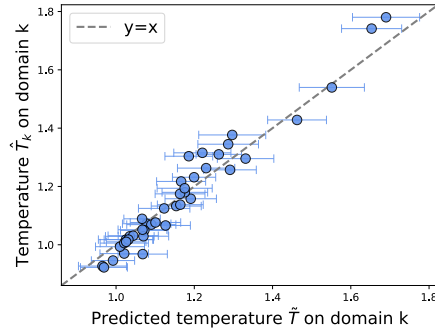


Figure 2: Compared the predicted temperature to the learned temperature \hat{T}_k on the k -th domain.

180 4 Experiments

181 In this section, we present experimental results evaluating our proposed method, demonstrating its
 182 effectiveness on both in-distribution and out-of-distribution calibration. We focus on three real-world
 183 datasets, including ImageNet-C [13]—a widely used robustness benchmark image classification
 184 dataset, WILDS-RxRx1 [18]—an image of cells (by fluorescent microscopy) dataset in the domain
 185 generalization benchmark, and GLDv2 [38]—a landmark recognition dataset in federated learning.
 186 Additional experimental results and implementation details can be found in Appendix B.

187 **Datasets.** We evaluate different calibration methods on three datasets, ImageNet-C, WILDS-RxRx1,
 188 and GLDv2. ImageNet-C contains 15 types of common corruptions where each corruption includes
 189 five severity levels. Each corruption with one severity is one domain, and there are 76 domains in
 190 total (including the standard ImageNet validation dataset). We partition the 76 domains into disjoint
 191 in-distribution domains and out-of-distribution by severity level or corruption type. WILDS-RxRx1 is
 192 a domain generalization dataset, and we treat each experimental domain as one domain. We adopt the
 193 default val/test split in Koh et al. [18]: use the four validation domains as in-distribution domains and
 194 the 14 test domains as the out-of-distribution domains. We also provide experimental results of other
 195 random splits in Appendix B. For GLDv2, each client corresponds to one domain, and there are 823
 196 domains in total. We randomly select 500 domains for training the model, and then use the remaining
 197 323 domains for evaluation denoted by validation domains. We further screen the validation domains
 198 by removing the domains with less than 300 data points. There are 44 domains after screening, and
 199 we use 30 domains as in-distribution domains and the remaining 14 domains as out-of-distribution
 200 domains. For all datasets, we randomly sample half of the data from in-distribution domains for
 201 calibrating models and use the remaining samples for InD ECE evaluation. We use all the samples
 202 from OOD domains for ECE evaluation.

203 **Models and training setup.** We consider multiple network architectures for evaluation, including
 204 ResNet-50 [12], ResNext-50 [39], DenseNet-121 [17], BiT-M-50 [19], Efficientnet-b1 [32], ViT-
 205 Small, and ViT-Base [7]. To evaluate on ImageNet-C, we directly evaluate models that are pre-trained
 206 on ImageNet [6]. For WILDS-RxRx1 and GLDv2, we use the ImageNet pre-trained models as
 207 initialization and apply SGD optimizer to training the models on training datasets.

208 **Evaluation metrics.** We use the Expected Calibration Error (ECE) as the main evaluation metric.
 209 We set the bin size as 100 for ImageNet-C, and set bin size as 20 for WILDS-RxRx1 and GLDv2.
 210 We evaluate ECE on both InD domains and OOD domains. Specifically, we evaluate the ECE of
 211 each InD/OOD domain. Meanwhile, we also evaluate the ECE of the pooled InD/OOD domains, i.e.,

Table 1: Per-domain ECE (%) comparison on three datasets. We evaluate the per-domain ECE on InD and OOD domains. We report the mean and standard error of per-domain ECE on one dataset. Lower ECE means better performance.

Datasets	Architectures	InD-domains			OOD-domains		
		MSP [14]	TS [11]	MD-TS	MSP [14]	TS [11]	MD-TS
ImageNet-C	ResNet-50	7.36±0.28	5.80±0.10	3.84±0.05	6.87±0.16	5.70±0.06	4.55±0.04
	Efficientnet-b1	6.78±0.07	6.12±0.15	3.99±0.07	6.54±0.06	4.87±0.05	4.05±0.03
	BiT-M-R50	6.93±0.27	6.99±0.25	3.86±0.06	6.32±0.16	6.50±0.16	4.30±0.04
	ViT-Base	4.77±0.16	4.34±0.12	3.76±0.07	4.09±0.06	4.01±0.05	3.86±0.04
WILDS-RxRx1	ResNet-50	26.22±0.38	9.83±0.57	2.85±0.17	26.22±0.38	13.78±0.43	5.25±0.11
	ResNext-50	25.30±0.76	9.39±0.58	3.13±0.19	20.71±0.30	11.80±0.37	5.07±0.09
	DenseNet-121	32.37±0.91	8.91±0.60	2.94±0.18	24.49±0.35	13.08±0.41	5.38±0.13
GLDv2	ResNet-50	12.56±0.08	11.61±0.09	9.90±0.06	11.36±0.15	10.75±0.14	9.76±0.12
	BiT-M-R50	14.86±0.12	11.31±0.07	9.78±0.06	13.91±0.21	9.83±0.11	9.16±0.10
	ViT-Small	12.44±0.11	11.12±0.07	9.75±0.05	11.00±0.18	9.65±0.11	9.01±0.10

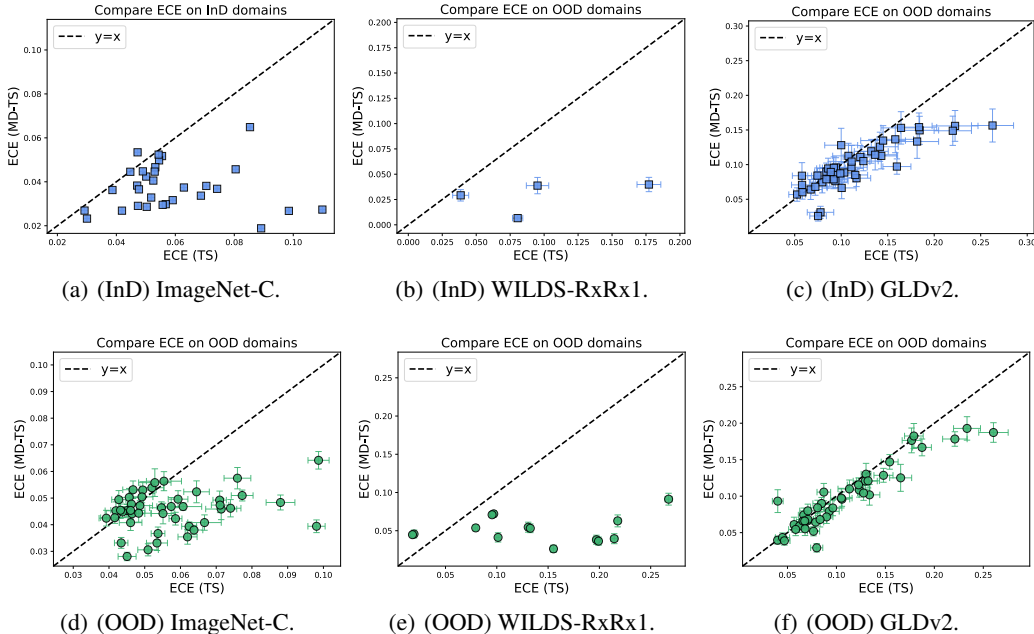


Figure 3: Per-domain ECE of MD-TS and TS on both in-distribution domains and out-of-distribution domains. Each plot is shown with ECE of TS (X -axis) and ECE of MD-TS (Y -axis). Top: per-domain ECE evaluated on InD domains. Bottom: per-domain ECE evaluated on OOD domains. Lower ECE is better.

212 the ECE evaluated on all samples from InD/OOD domains. We use unseen samples from the InD
 213 domain to measure the per-domain ECE. We also measure the averaged per-domain ECE results (i.e.,
 214 per-domain ECE averaged across domains).

215 4.1 Main results

216 We summarize the ECE results of different methods on three datasets in Table 1 and Figure 3. We
 217 use TS to denote temperature scaling [11], and use MSP to denote applying the maximum softmax
 218 probability [14] of the model output (i.e., without calibration). In Table 1, we use the ImageNet
 219 validation dataset and ImageNet-C datasets with severity level $s \in \{1, 5\}$ as the InD domains and use
 220 the remaining datasets as OOD domains. We present the averaged per-domain ECE results in Table 1,
 221 and visualize the ECE of each domain in Figure 3. As shown in Table 1 and Figure 3(a)-3(c), we
 222 find that our proposed approach achieves much better InD per-domain calibration compared with

Table 2: Model performance prediction comparison results of different methods on three datasets. Lower MAE indicates better performance.

Datasets	Architectures	InD-domains MAE			OOD-domains MAE		
		MSP [14]	TS [11]	MD-TS	MSP [14]	TS [11]	MD-TS
ImageNet-C	ResNet-50	5.88	4.74	1.28	5.15	3.96	1.70
	BiT-M-R50	6.08	6.16	1.33	4.97	5.23	1.66
WILDS-RxRx1	ResNet-50	33.65	9.61	1.61	26.20	13.66	4.76
	ResNext-50	25.32	8.55	1.39	20.72	12.88	4.78
GLDv2	ResNet-50	9.60	9.17	7.11	9.72	9.40	8.08
	BiT-M-R50	12.67	7.18	4.64	12.30	7.34	6.37

223 baselines. Also, TS does not significantly improve over MSP on ImageNet-C InD domains in Table 1,
 224 but our proposed method largely improve the ECE compared with MSP and TS. For instance, the
 225 ECE results of MSP and TS on Efficientnet-b1 are 6.93 and 6.99, and our method achieves 3.84.
 226 Intuitively, when there are a diverse set of domains in the calibration dataset, a single temperature
 227 cannot provide well-calibrated confidences. In contrast, our proposed method is able to produce
 228 much better InD confidence estimates by leveraging the multi-domain structure of the data.

229 Next we study the performance of different methods on out-of-distribution domains. From Table 1,
 230 we find that MD-TS achieves the best performance on OOD domains across all the settings. On
 231 ImageNet-C with BiT-M-R50, MD-TS improves the ECE from 6.54 (MSP) to 4.05, while the
 232 performance of TS is similar to MSP. Moreover, MD-TS significantly outperforms MSP and TS on
 233 WILDS-RxRx1, where MD-TS improves over TS by around 5.00 measured in ECE. Figure 3(d)-3(f)
 234 display the per-domain ECE performance on out-of-distribution domains. MD-TS improves over TS
 235 on more than half of the domains in all three datasets. For the remaining domains, MD-TS performs
 236 slightly worse than TS. Furthermore, on those domains that TS performs poorly (ECE > 8), MD-TS
 237 largely improves over TS by large margins.

238 4.2 Predicting generalization

239 Suppose a model can produce calibrated confidences on unseen samples, in which case we could
 240 leverage the calibrated confidence to predict the model performance. Specifically, based on the
 241 definition of ECE in Eq. (1), when the model is well-calibrated, the average of the calibrated
 242 confidence is close to the average accuracy, i.e., $\text{Conf}(\mathcal{D}) \approx \text{Acc}(\mathcal{D})$.² Meanwhile, predicting model
 243 performance accurately is an essential ingredient in developing reliable machine learning systems,
 244 especially under distributional shifts [10]. As shown in Table 1, we find that our proposed method
 245 produces well-calibrated confidence values on both InD and OOD domains. We now measure its
 246 performance on predicting model performance and compare with existing methods. We measure
 247 the performance using mean absolute error (MAE), $\text{MAE} = (1/K) \cdot \sum_{k=1}^K |\text{Conf}(\mathcal{D}_k) - \text{Acc}(\mathcal{D}_k)|$
 248 where S_k is the dataset from the k -th domain.

249 We show the predicting model accuracy results in Table 2. MD-TS significantly improves over existing
 250 methods on predicting model performance across all three datasets. For example, on ImageNet-C,
 251 calibrated confidence of MD-TS produces fairly accurate predictions on both InD and OOD domains
 252 (less than 2% measured in MAE), which largely outperforms MSP and TS. In Figure 4, we compared
 253 the prediction performance of TS and MD-TS on every OOD domain. We find that MD-TS achieves
 254 better prediction performance compared to TS on most of the domains. Refer to Appendix B.1 for
 255 more results in which other architectures are tested.

256 4.3 MD-TS ablations

257 To learn a calibration model that performs well per-domain, we apply linear regression on feature
 258 representations $\Phi(x_k)$ such that $\langle \Phi(x_k), \theta \rangle \approx \hat{T}_k$, where x_k is from domain k and \hat{T}_k is the tem-
 259 perature parameter for domain k . We investigate other methods for learning the map from feature

² $\text{Conf}(\mathcal{D})$ denotes the average (calibrated) confidence on dataset \mathcal{D} , and $\text{Acc}(\mathcal{D})$ denotes the average accuracy on dataset \mathcal{D} .

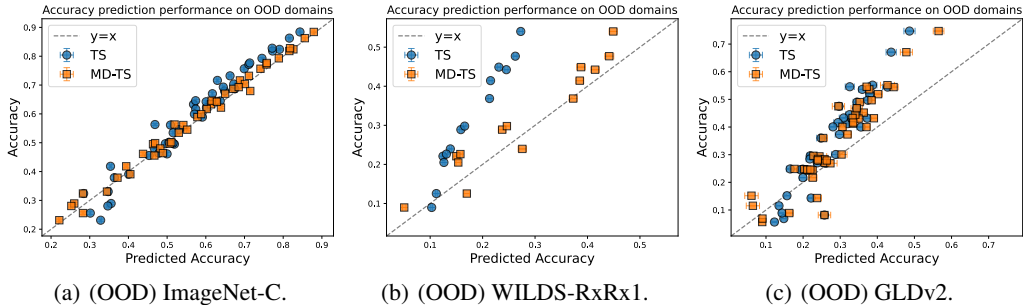


Figure 4: Predicting accuracy performance of MD-TS and TS on both out-of-distribution domains. Each plot is shown with predicted accuracy (X -axis) and accuracy (Y -axis). Each points corresponds to one domain. The network architecture is ResNet-50 for three datasets. Point closer to the $Y = X$ dashed line means better prediction performance.

Table 3: Per-domain ECE (%) results of MD-TS ablations on WILDS-RxRx1. We evaluate the per-domain ECE on InD and OOD domains, and report the mean and standard error of per-domain ECE. Lower ECE means better performance.

Architectures	InD-domains					OOD-domains				
	OLS	Ridge	Huber	KRR	KNN	OLS	Ridge	Huber	KRR	KNN
ResNet-50	2.85	2.88	2.90	2.85	3.00	5.25	5.26	5.29	4.99	5.44
ResNext-50	3.13	3.14	3.11	3.07	3.03	5.07	5.06	5.02	4.94	5.36
DenseNet-121	2.94	3.03	2.92	2.90	3.04	5.38	5.42	5.36	5.20	5.47

260 representations to temperatures in a regression framework. Specifically, beside the ordinary least
 261 squares (OLS) used in Algorithm 1, we consider ridge regression (Ridge), robust regression with
 262 Huber loss (Huber), kernel ridge regression (KRR), and K -nearest neighbors regression (KNN). The
 263 implementations are mainly based on `scikit-learn` [28]. We use grid search (on InD domains) to
 264 select hyperparameters for Ridge, Huber, KRR, and KNN.

265 We summarize the comparative results for different regression algorithms in Table 3. Compared to
 266 OLS, other regression algorithms do not achieve significant improvement. Specifically, KRR achieves
 267 slightly better performance on OOD domains, while other algorithms have similar performance
 268 compared to OLS. Moreover, there are no hyperparameter in OLS, which makes it more practical in
 269 real-world problems. Meanwhile, the results suggest that our proposed MD-TS is stable to the choice
 270 of specific regression algorithms.

271 5 Theoretical analysis

272 In this section, we provide theoretical analysis to support our understanding of our proposed algorithm
 273 in the presence of distribution shifts. We use $h_k^*(\cdot) = h(\cdot; f, \beta_k^*) : \mathcal{X} \rightarrow [0, 1]$ to denote the best
 274 calibration map for the base model f on the k -th domain; this map *minimizes* the expected calibration
 275 error (ECE) $\mathbb{E}[|p - \mathbb{P}(\hat{y} = y | \hat{\pi} = p)|]$ over distribution \mathcal{P}_k . We also call h_k^* a hypothesis in the
 276 hypothesis class \mathcal{H} . Next, given the fixed base model f , we aim to learn $\hat{h}(\cdot) = h(\cdot; f, \hat{\beta})$ such
 277 that $\varepsilon(\hat{h}, \mathcal{P}_{k,X}) = \mathbb{E}_{X \sim \mathcal{P}_{k,X}}[|h_k^*(X) - \hat{h}(X)|]$ is small for *every* domain k , where $\varepsilon_k(\hat{h})$ denotes the
 278 risk of \hat{h} w.r.t. the the best calibration map h_k^* under domain \mathcal{P}_k . In addition, we are interested in
 279 generalizing to new domains: suppose there is an unseen OOD domain $\tilde{\mathcal{P}}$ and its marginal feature
 280 distribution is different from existing domains, i.e., $\tilde{\mathcal{P}}_X \neq \mathcal{P}_{k,X}$ for $k \in [K]$.

281 Our goal is to understand the conditions under which \hat{h} can have similar calibration on OOD domains
 282 as the InD domains. For example, if the OOD domain is similar to the mixture distribution of InD
 283 domains, we would expect \hat{h} performs similarly on InD and OOD domains. To quantify the distance
 284 between two distributions, we first introduce the \mathcal{H} -divergence [3] to measure the distance between
 285 two distributions:

286 **Definition 5.1** (\mathcal{H} -divergence). *Given an input space \mathcal{X} and two probability distributions \mathcal{P}_X and*
 287 *\mathcal{P}'_X on \mathcal{X} , let \mathcal{H} be a hypothesis class on \mathcal{X} , and denote by \mathcal{A} the collection of subsets of \mathcal{X} which*
 288 *are the support of hypothesis $h \in \mathcal{H}$, i.e., $\mathcal{A}_{\mathcal{H}} = \{h^{-1}(1) \mid h \in \mathcal{H}\}$. The distance between \mathcal{P}_X and*
 289 *\mathcal{P}'_X is defined as*

$$d_{\mathcal{H}}(\mathcal{P}_X, \mathcal{P}'_X) = \sup_{A \in \mathcal{A}_{\mathcal{H}}} |\Pr_{\mathcal{P}_X}(A) - \Pr_{\mathcal{P}'_X}(A)|.$$

290 The \mathcal{H} -divergence reduces to the standard total variation (TV) distance when \mathcal{H} contains all mea-
 291 surable functions on \mathcal{X} , which implies that the \mathcal{H} -divergence is upper bounded by the TV-distance,
 292 i.e., $d_{\mathcal{H}}(\mathcal{P}_X, \mathcal{P}'_X) \leq d_{\text{TV}}(\mathcal{P}_X, \mathcal{P}'_X)$. On the other hand, when the hypothesis class \mathcal{H} has a finite VC
 293 dimension or pseudo-dimension, the \mathcal{H} -divergence can be estimated using finite samples from \mathcal{P}_X
 294 and \mathcal{P}'_X [3]. Next, we define the mixture distribution of the K in-distribution domains $\mathcal{P}_{K,X}^\alpha$ on input
 295 space \mathcal{X} as follows:

$$\mathcal{P}_{K,X}^\alpha = \sum_{k=1}^K \alpha_k \mathcal{P}_{k,X}, \quad \text{where } \sum_{k=1}^K \alpha_k = 1 \text{ and } \alpha_k \geq 0.$$

296 Given multiple domains $\{\mathcal{P}_1, \dots, \mathcal{P}_K\}$, we can optimize the combination parameters α such that
 297 $\mathcal{P}_{K,X}^\alpha$ minimizes the \mathcal{H} -divergence between $\mathcal{P}_{K,X}^\alpha$ and $\tilde{\mathcal{P}}_X$. More specifically, we define $\hat{\alpha}$ as

$$\hat{\alpha} = \underset{\alpha \in \Delta}{\operatorname{argmin}} \left\{ \frac{1}{2} d_{\mathcal{H}}(\mathcal{P}_{K,X}^\alpha, \tilde{\mathcal{P}}_X) + \lambda(\mathcal{P}_{K,X}^\alpha, \tilde{\mathcal{P}}_X) \right\}, \quad \lambda(\mathcal{P}_{K,X}^\alpha, \tilde{\mathcal{P}}_X) = \varepsilon(h^*, \mathcal{P}_{K,X}^\alpha) + \varepsilon(h^*, \tilde{\mathcal{P}}_X), \quad (3)$$

298 where $h^* := \operatorname{argmin}_{h \in \mathcal{H}} \{\varepsilon(h, \mathcal{P}_{K,X}^\alpha) + \varepsilon(h, \tilde{\mathcal{P}}_X)\}$ and $\tilde{\mathcal{H}}$ is defined as $\tilde{\mathcal{H}} := \{\operatorname{sign}(|h(x) - h'(x)| -$
 299 $t) \mid h, h' \in \mathcal{H}, 0 \leq t \leq 1\}$. We now give an upper bound on the risk on the unseen OOD domain. This
 300 result follows very closely those of Blitzer et al. [4], Zhao et al. [42], instantiated in our calibration
 301 setup. Details can be found in Appendix C.

302 **Theorem 5.2.** *Let \mathcal{H} be a hypothesis class that contains functions $h : \mathcal{X} \rightarrow [0, 1]$ with pseudo-*
 303 *dimension $\operatorname{Pdim}(\mathcal{H}) = d$. Let $\{\mathcal{D}_{k,X}\}_{k=1}^K$ denote the empirical distributions generated from*
 304 *$\{\mathcal{P}_{k,X}\}_{k=1}^K$, where $\mathcal{D}_{k,X}$ contains n i.i.d. samples from the marginal feature distribution $\mathcal{P}_{k,X}$ of*
 305 *domain k . Then for $\delta \in (0, 1)$, with probability at least $1 - \delta$, we have*

$$\varepsilon(\hat{h}, \tilde{\mathcal{P}}_X) \leq \sum_{k=1}^K \hat{\alpha}_k \cdot \hat{\varepsilon}(\hat{h}, \mathcal{D}_{k,X}) + \frac{1}{2} d_{\mathcal{H}}(\mathcal{P}_{K,X}^\alpha, \tilde{\mathcal{P}}_X) + \lambda(\mathcal{P}_{K,X}^\alpha, \tilde{\mathcal{P}}_X) + \tilde{O}\left(\frac{\operatorname{Pdim}(\mathcal{H})}{\sqrt{nK}}\right), \quad (4)$$

306 where $\hat{\alpha}$ and $\lambda(\mathcal{P}_{K,X}^\alpha, \tilde{\mathcal{P}}_X)$ are defined in Eq. (3), $\tilde{\mathcal{P}}_X$ denotes the marginal distribution of the
 307 OOD domain, $\operatorname{Pdim}(\mathcal{H})$ is the pseudo-dimension of the hypothesis class \mathcal{H} , and $\hat{\varepsilon}(\hat{h}, \mathcal{D}_{k,X})$ is the
 308 empirical risk of the hypothesis \hat{h} on $\mathcal{D}_{k,X}$.

309 This result means that if we can learn a hypothesis \hat{h} that achieves small empirical risk $\hat{\varepsilon}(\hat{h}, \mathcal{D}_{k,X})$
 310 on every domain, then \hat{h} is able to achieve good performance on the OOD domain if distribution of
 311 the OOD domain is similar to the mixture distribution of InD domains measured by \mathcal{H} -divergence.
 312 In this case, if the learned calibration map \hat{h} is well-calibrated on every domain \mathcal{P}_k , then \hat{h} is likely
 313 to provide calibrated confidence for the OOD domain $\tilde{\mathcal{P}}$. Recall from Section 4, we proposed an
 314 algorithm that performs well across InD domains. The upper bound in Eq. (4) provides insight into
 315 understanding why this algorithm is effective.

316 6 Discussion

317 We have developed an algorithm for robust calibration that exploits multi-domain structure in datasets.
 318 Experiments on real-world domains indicate that multi-domain calibration is an effective way to
 319 improve the robustness of calibration under distribution shifts. One interesting direction for future
 320 work would be to extend our algorithm to a scenario where no domain information is available. We
 321 hope the multi-domain calibration perspective in this paper can motivate further work to close the
 322 gap between in-distribution and out-of-distribution calibration.

323 References

- 324 [1] Anastasios N. Angelopoulos, Stephen Bates, Emmanuel J. Candès, Michael I. Jordan, and Lihua
325 Lei. Learn then test: Calibrating predictive algorithms to achieve risk control. *arXiv preprint*,
326 2021. arXiv:2110.01052.
- 327 [2] Stephen Bates, Anastasios Angelopoulos, Lihua Lei, Jitendra Malik, and Michael I. Jordan.
328 Distribution-free, risk-controlling prediction sets. *Journal of the ACM*, 68(6), September 2021.
329 doi: 10.1145/3478535.
- 330 [3] Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jen-
331 nifer Wortman Vaughan. A theory of learning from different domains. *Machine learning*, 79(1):
332 151–175, 2010.
- 333 [4] John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman. Learning
334 bounds for domain adaptation. *Advances in neural information processing systems*, 20, 2007.
- 335 [5] Maxime Cauchois, Suyash Gupta, Alnur Ali, and John C Duchi. Robust validation: Confident
336 predictions even when distributions shift. *arXiv preprint arXiv:2008.04267*, 2020.
- 337 [6] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-
338 scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern
339 recognition*, pages 248–255. Ieee, 2009.
- 340 [7] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
341 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al.
342 An image is worth 16x16 words: Transformers for image recognition at scale. In *International
343 Conference on Learning Representations*, 2020.
- 344 [8] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model
345 uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059.
346 PMLR, 2016.
- 347 [9] Yarin Gal, Jiri Hron, and Alex Kendall. Concrete dropout. *Advances in neural information
348 processing systems*, 30, 2017.
- 349 [10] Devin Guillory, Vaishal Shankar, Sayna Ebrahimi, Trevor Darrell, and Ludwig Schmidt. Pre-
350 dicting with confidence on unseen distributions. In *Proceedings of the IEEE/CVF International
351 Conference on Computer Vision*, pages 1134–1144, 2021.
- 352 [11] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural
353 networks. In *International Conference on Machine Learning*, pages 1321–1330. PMLR, 2017.
- 354 [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
355 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
356 pages 770–778, 2016.
- 357 [13] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common
358 corruptions and perturbations. *arXiv preprint arXiv:1903.12261*, 2019.
- 359 [14] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution
360 examples in neural networks. *arXiv preprint arXiv:1610.02136*, 2016.
- 361 [15] Dan Hendrycks, Mantas Mazeika, Saurav Kadavath, and Dawn Song. Using self-supervised
362 learning can improve model robustness and uncertainty. *Advances in Neural Information
363 Processing Systems*, 32, 2019.
- 364 [16] Dan Hendrycks, Norman Mu, Ekin D Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshmi-
365 narayanan. Augmix: A simple data processing method to improve robustness and uncertainty.
366 *arXiv preprint arXiv:1912.02781*, 2019.
- 367 [17] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected
368 convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern
369 recognition*, pages 4700–4708, 2017.

- 370 [18] Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay
371 Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, et al.
372 Wilds: A benchmark of in-the-wild distribution shifts. In *International Conference on Machine*
373 *Learning*, pages 5637–5664. PMLR, 2021.
- 374 [19] Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain
375 Gelly, and Neil Houlsby. Big transfer (bit): General visual representation learning. In *European*
376 *conference on computer vision*, pages 491–507. Springer, 2020.
- 377 [20] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable
378 predictive uncertainty estimation using deep ensembles. *Advances in neural information*
379 *processing systems*, 30, 2017.
- 380 [21] Donghwan Lee, Xinmeng Huang, Hamed Hassani, and Edgar Dobriban. T-cal: An optimal test
381 for the calibration of predictive models. *arXiv preprint arXiv:2203.01850*, 2022.
- 382 [22] Matthias Minderer, Josip Djolonga, Rob Romijnders, Frances Hubis, Xiaohua Zhai, Neil
383 Houlsby, Dustin Tran, and Mario Lucic. Revisiting the calibration of modern neural networks.
384 *Advances in Neural Information Processing Systems*, 34, 2021.
- 385 [23] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of machine learning*.
386 2018.
- 387 [24] Mahdi Pakdaman Naeini, Gregory Cooper, and Milos Hauskrecht. Obtaining well calibrated
388 probabilities using bayesian binning. In *Twenty-Ninth AAAI Conference on Artificial Intelligence*,
389 2015.
- 390 [25] Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, David Sculley, Sebastian Nowozin, Joshua
391 Dillon, Balaji Lakshminarayanan, and Jasper Snoek. Can you trust your model’s uncertainty?
392 evaluating predictive uncertainty under dataset shift. *Advances in neural information processing*
393 *systems*, 32, 2019.
- 394 [26] Sangdon Park, Osbert Bastani, James Weimer, and Insup Lee. Calibrated prediction with
395 covariate shift via unsupervised domain adaptation. In *International Conference on Artificial*
396 *Intelligence and Statistics*, pages 3219–3229. PMLR, 2020.
- 397 [27] Sangdon Park, Edgar Dobriban, Insup Lee, and Osbert Bastani. Pac prediction sets under
398 covariate shift. *arXiv preprint arXiv:2106.09848*, 2021.
- 399 [28] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel,
400 P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher,
401 M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine*
402 *Learning Research*, 12:2825–2830, 2011.
- 403 [29] John Platt et al. Probabilistic outputs for support vector machines and comparisons to regularized
404 likelihood methods. *Advances in large margin classifiers*, 10(3):61–74, 1999.
- 405 [30] Aleksandr Podkopaev and Aaditya Ramdas. Distribution-free uncertainty quantification for
406 classification under label shift. In *Uncertainty in Artificial Intelligence*, pages 844–853. PMLR,
407 2021.
- 408 [31] Joaquin Quiñero-Candela, Masashi Sugiyama, Anton Schwaighofer, and Neil D Lawrence.
409 *Dataset shift in machine learning*. Mit Press, 2008.
- 410 [32] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural
411 networks. In *International conference on machine learning*, pages 6105–6114. PMLR, 2019.
- 412 [33] Sunil Thulasidasan, Gopinath Chennupati, Jeff A Bilmes, Tanmoy Bhattacharya, and Sarah
413 Michalak. On mixup training: Improved calibration and predictive uncertainty for deep neural
414 networks. *Advances in Neural Information Processing Systems*, 32, 2019.
- 415 [34] Ryan J Tibshirani, Rina Foygel Barber, Emmanuel Candes, and Aaditya Ramdas. Conformal
416 prediction under covariate shift. *Advances in neural information processing systems*, 32, 2019.

- 417 [35] Vladimir Vovk, Alexander Gammerman, and Glenn Shafer. *Algorithmic learning in a random*
418 *world*. Springer Science & Business Media, 2005.
- 419 [36] Yoav Wald, Amir Feder, Daniel Greenfeld, and Uri Shalit. On calibration and out-of-domain
420 generalization. *Advances in Neural Information Processing Systems*, 34, 2021.
- 421 [37] Ximei Wang, Mingsheng Long, Jianmin Wang, and Michael Jordan. Transferable calibration
422 with lower bias and variance in domain adaptation. *Advances in Neural Information Processing*
423 *Systems*, 33:19212–19223, 2020.
- 424 [38] Tobias Weyand, Andre Araujo, Bingyi Cao, and Jack Sim. Google landmarks dataset v2-a large-
425 scale benchmark for instance-level recognition and retrieval. In *Proceedings of the IEEE/CVF*
426 *conference on computer vision and pattern recognition*, pages 2575–2584, 2020.
- 427 [39] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual
428 transformations for deep neural networks. In *Proceedings of the IEEE conference on computer*
429 *vision and pattern recognition*, pages 1492–1500, 2017.
- 430 [40] Bianca Zadrozny and Charles Elkan. Obtaining calibrated probability estimates from decision
431 trees and naive bayesian classifiers. In *Icml*, volume 1, pages 609–616. Citeseer, 2001.
- 432 [41] Bianca Zadrozny and Charles Elkan. Transforming classifier scores into accurate multiclass
433 probability estimates. In *Proceedings of the eighth ACM SIGKDD international conference on*
434 *Knowledge discovery and data mining*, pages 694–699, 2002.
- 435 [42] Han Zhao, Shanghang Zhang, Guanhang Wu, José MF Moura, Joao P Costeira, and Geoffrey J
436 Gordon. Adversarial multiple source domain adaptation. *Advances in neural information*
437 *processing systems*, 31, 2018.

438 **Checklist**

- 439 1. For all authors...
- 440 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
441 contributions and scope? [Yes]
- 442 (b) Did you describe the limitations of your work? [Yes]
- 443 (c) Did you discuss any potential negative societal impacts of your work? [No]
- 444 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
445 them? [Yes]
- 446 2. If you are including theoretical results...
- 447 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See Theo-
448 rem 5.2.
- 449 (b) Did you include complete proofs of all theoretical results? [Yes] See Appendix C.
- 450 3. If you ran experiments...
- 451 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
452 mental results (either in the supplemental material or as a URL)? [Yes] See supplemen-
453 tal material.
- 454 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
455 were chosen)? [Yes] See Section 4 and Appendix A.
- 456 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
457 ments multiple times)? [Yes] See Section 4.
- 458 (d) Did you include the total amount of compute and the type of resources used (e.g., type
459 of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix A.
- 460 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 461 (a) If your work uses existing assets, did you cite the creators? [Yes] See Section 4.
- 462 (b) Did you mention the license of the assets? [N/A]
- 463 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 464
- 465 (d) Did you discuss whether and how consent was obtained from people whose data you're
466 using/curating? [N/A]
- 467 (e) Did you discuss whether the data you are using/curating contains personally identifiable
468 information or offensive content? [N/A]
- 469 5. If you used crowdsourcing or conducted research with human subjects...
- 470 (a) Did you include the full text of instructions given to participants and screenshots, if
471 applicable? [N/A]
- 472 (b) Did you describe any potential participant risks, with links to Institutional Review
473 Board (IRB) approvals, if applicable? [N/A]
- 474 (c) Did you include the estimated hourly wage paid to participants and the total amount
475 spent on participant compensation? [N/A]