
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 [37]. 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
20 that is valid for *every* hospital—not just on average across all hospitals. Going even further, our
21 uncertainty quantification should be informative when a new 21st hospital goes online, even if the
22 distribution shifts from those already encountered. *As another example, a centralized model is trained
23 on training data from existing clients in federated learning [8]. It is important for the central server
24 to provide uncertainty quantification for every client. Similar to the fMRI example, the centralized
25 model should still produce valid uncertainty quantification for unseen new clients. Another example
26 is applying animal recognition models on images in wildlife monitoring [4], where one set of camera
27 traps corresponds to one domain, and the model will be deployed under distribution shift, i.e., new
28 camera traps.* In this work, we study calibration in the multi-domain setting. We find that by requiring
29 accurate calibration across all observed domains, our method provides more accurate uncertainty
30 quantification on unseen domains.

31 Calibration is a core topic in learning [35, 30, 11, 26, 14, 3], but most techniques are targeted at
32 settings with no distribution shift. To see this, we consider a simple experiment on the ImageNet-
33 C [16] dataset, which consists of 76 domains. Here, each domain corresponds to one type of data
34 corruption applied with a certain severity. We apply the temperature scaling technique [14] on the
35 pooled data from all domains. In Figure 1(a) and 1(b), we display the reliability diagrams for the
36 pooled data and for one individual domain. We find that even under a relatively mild distribution

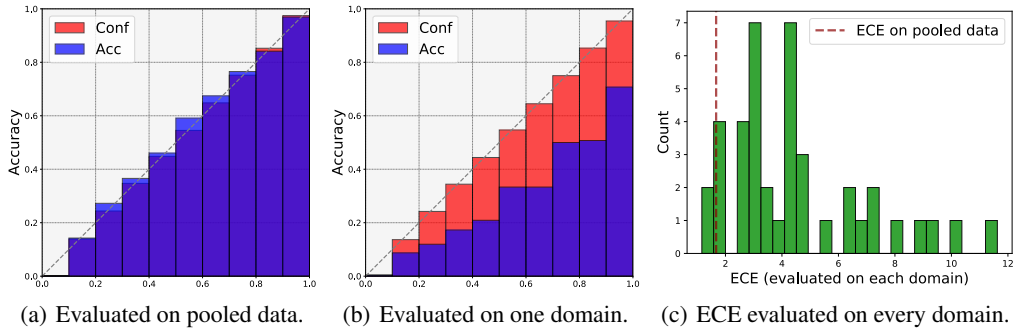


Figure 1: Reliability diagrams and expected calibration error histograms for temperature scaling with a ResNet-50 on ImageNet-C. We use temperature scaling to obtain adjusted confidences for the ResNet-50 model. **(a)** Reliability diagram evaluated on the pooled data of ImageNet-C. **(b)** Reliability 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 shift—i.e., subpopulation shift from the mixture of all domains to the single domain—temperature
 38 scaling does not produce calibrated confidence estimates on the stand-alone domain. This behavior is
 39 pervasive; in Figure 1(c), we see that the calibration on individual domains is much worse than the
 40 the reliability diagram from the pooled data would suggest.

41 To address this issue, we develop a new algorithm, multi-domain temperature scaling, that leverages
 42 multi-domain structure in the data. Our algorithm takes a base model and learns a calibration function
 43 that maps each input to a different temperature parameter that is used for adjusting confidence in the
 44 base model. Empirically, we find our algorithm significantly outperforms temperature scaling on
 45 three real-world multi-domain datasets. In particular, in contrast to temperature scaling, our proposed
 46 algorithm is able to provide well-calibrated confidence on every domain. Moreover, our algorithm
 47 largely improves robustness of calibration under distribution shifts. This is expected, because if the
 48 calibration method performs well on every domain, it is likely to have learned some structure that
 49 generalizes to unseen domains. Theoretically, we analyze the multi-domain calibration problem in
 50 the regression setting, providing guidance about the conditions under which robust calibration is
 51 possible.

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

60 Related Work

61 **Calibration methods.** There is a large literature on calibrating the well-trained machine learning
 62 models, including histogram binning [47], isotonic regression [48], conformal prediction [41], Platt
 63 scaling [35], and temperature scaling [14]. These calibration methods apply a validation set and post-
 64 process the model outputs. As shown in Guo et al. [14], temperature scaling, a simple method that uses
 65 a single (temperature) parameter for rescaling the logits, performs surprisingly well on calibrating
 66 confidences for deep neural networks. We focus on this approach in our work. More broadly,
 67 there has been much recent work develop methods to improve calibration for deep learning models,
 68 including augmentation-based training [39, 19], [calibration for neural machine translation](#) [25], [neural](#)
 69 [stochastic differential equation](#) [24], self-supervised learning [18], ensembling [26], and Bayesian
 70 neural networks [11, 12], as well as statistical guarantees for calibration with black-box models [1].

71 **Calibration under distribution shifts.** Ovadia et al. [31] conduct an empirical study on model
 72 calibration under distribution shifts and find that models are much less calibrated under distribution
 73 shifts. Minderer et al. [28] revisit calibration of recent state-of-the-art image classification models
 74 under distribution shifts and study the relationship between calibration and accuracy. Wald et al.
 75 [42] study model calibration and out-of-distribution generalization. Other works consider providing
 76 uncertainty estimates under structured distribution shifts, such as covariate shift [40, 33], label

77 shift [36], and f -divergence balls [7]. Another line of work studies calibration in the domain
78 adaptation setting [43, 32], which require unlabeled samples from the target domain.

79 2 Problem setup

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

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

99 **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
100 (empirical) expected calibration error (ECE) with M bins evaluated on \mathcal{D} is defined as

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

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

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

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

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

103 The empirical ECE defined in Eq. (1) approximates the expected calibration error (ECE) $\mathbb{E}[|p - \mathbb{P}(\hat{y} =$
104 $y | \hat{\pi} = p)|]$ with bin size equal to M [30, 14]; see [27] for statistical results about about the empirical
105 ECE as an estimator. The perfect calibrated map corresponds to the case when $\mathbb{P}(\hat{y} = y | \hat{\pi} = p) = p$
106 holds for all $p \in [0, 1]$.

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

118 **Definition 2.2** (Multi-domain ECE). For a dataset $\mathcal{D}_K^{\text{pool}} = \mathcal{D}_1 \cup \dots \cup \mathcal{D}_K$ consisting of samples
119 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)
120 multi-domain expected calibration error (Multi-domain ECE) with M bins evaluated on $\mathcal{D}_K^{\text{pool}}$ is
121 defined as $\text{MDECE}(\mathcal{D}_K^{\text{pool}}) = \frac{1}{K} \sum_{k=1}^K \text{ECE}(\mathcal{D}_k)$.

122 **Remark 2.3.** *In Definition 2.2, we weight each domain equally to balance across domains, which*
 123 *could better reflect how the calibration method performs on each individual domain. Furthermore,*
 124 *in our experiments, we also visualize the ECE measured on each domain to provide additional*
 125 *information on model performance on every domain.*

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

132 **Temperature scaling.** Next, we review a simple and effective calibration method, named temperature
 133 scaling (TS) [35, 14], that is widely used in single-domain model calibration. Temperature scaling
 134 applies a single parameter $T > 0$ and produces the confidence prediction for the base model f as

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

135 where $[\text{Softmax}(z)]_j = \exp([z]_j) / \sum_{i=1}^J \exp([z]_i)$. The parameter T is the so-called *temperature*,
 136 with larger temperature yielding more diffuse probability estimates. To learn the temperature
 137 parameter T from dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$, Guo et al. [14] propose to find T by solving the
 138 following convex optimization problem,

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

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

142 3 Multi-domain temperature scaling

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

149 At a high level, we propose to learn a calibration model that maps samples from the input space \mathcal{X} to
 150 the temperature space \mathbb{R}_+ . To start with, we learn the temperature parameter \hat{T}_k for the base model
 151 on every domain k by applying temperature scaling on \mathcal{D}_k . Next, we apply the base deep model to
 152 compute feature embeddings of samples from different domains,¹ and label feature embeddings from
 153 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
 154 contains feature embeddings and temperature labels from one domain, i.e., $\hat{\mathcal{D}}_k = \{(\Psi(x_{i,k}), \hat{T}_k)\}_{i=1}^{n_k}$.
 155 Finally, we apply linear regression on these labeled datasets. In detail, our algorithm is as follows:

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

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

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

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

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

165 We denote our proposed method by MD-TS (**Mult-Domain Temperature Scaling**). A presentation of
 166 the algorithm in pseudocode can be found in Algorithm 1, Appendix A.

167 We pause to consider the basic concept in more detail. The goal of our proposed algorithm is to
 168 predict the best temperature for samples from different several domains. In an ideal setting where
 169 the learned linear model $\hat{\theta}$ results in good calibration on *every* InD domain, we can expect that $\hat{\theta}$
 170 will continue to yield good calibration on the OOD domain $\tilde{\mathcal{P}}$ when $\tilde{\mathcal{P}}$ is close to one or several InD
 171 domains. For example, $\tilde{\mathcal{P}}$ will work well if $\tilde{\mathcal{P}}$ is a mixture of the K domains, i.e., $\tilde{\mathcal{P}} = \sum_{k=1}^K \alpha_k \mathcal{P}_k$
 172 and $\alpha \in \Delta^{K-1}$. Regarding the algorithmic design, linear regression is one of the simplest models
 173 for solving the regression problem. It is computationally fast to learn such linear models as well as
 174 make predictions on new samples, making it attractive. We test alternative, more flexible, regression
 175 algorithms in Section 4 but do not observe significant gains over linear regression.

176 To illustrate how our proposed algorithm MD-TS performs differently from standard TS, we return to the
 177 ImageNet-C dataset. We compare the predicted temperature of our algorithm on new samples from domain
 178 k with the temperature that results from running TS on domain k alone. The results are summarized in
 179 Figure 2, where each circle corresponds to the mean predicted temperature on one InD domain. For each
 180 domain, we also visualize the standard deviation of the predicted temperatures for samples from that domain
 181 (the horizontal bar around each point). We find that our algorithm predicts the temperature quite well. Note that
 182 it does not have access to the domain index information of the fresh samples. By contrast, TS always uses the
 183 same temperature, regardless of the input point.
 184
 185
 186
 187
 188
 189
 190

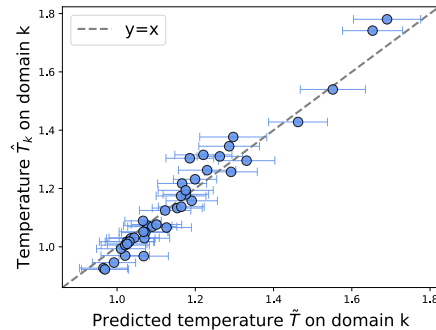


Figure 2: Compare the predicted temperature to the learned temperature \hat{T}_k on the k -th domain.

191 4 Experiments

192 In this section, we present experimental results evaluating our proposed method, demonstrating its
 193 effectiveness on both in-distribution and out-of-distribution calibration. We focus on three real-world
 194 datasets, including ImageNet-C [16]—a widely used robustness benchmark image classification
 195 dataset, WILDS-RxRx1 [22]—an image of cells (by fluorescent microscopy) dataset in the domain
 196 generalization benchmark, and GLDv2 [44]—a landmark recognition dataset in federated learning.
 197 Additional experimental results and implementation details can be found in Appendix B.

198 **Datasets.** We evaluate different calibration methods on three datasets, ImageNet-C, WILDS-RxRx1,
 199 and GLDv2. ImageNet-C contains 15 types of common corruptions where each corruption includes
 200 five severity levels. Each corruption with one severity is one domain, and there are 76 domains in
 201 total (including the standard ImageNet validation dataset). We partition the 76 domains into disjoint
 202 in-distribution domains and out-of-distribution by severity level or corruption type. WILDS-RxRx1 is
 203 a domain generalization dataset, and we treat each experimental domain as one domain. We adopt the
 204 default val/test split in Koh et al. [22]: use the four validation domains as in-distribution domains and
 205 the 14 test domains as the out-of-distribution domains. We also provide experimental results of other
 206 random splits in Appendix B. For GLDv2, each client corresponds to one domain, and there are 823
 207 domains in total. We randomly select 500 domains for training the model, and then use the remaining
 208 323 domains for evaluation denoted by validation domains. We further screen the validation domains
 209 by removing the domains with less than 300 data points. There are 44 domains after screening, and
 210 we use 30 domains as in-distribution domains and the remaining 14 domains as out-of-distribution
 211 domains. For all datasets, we randomly sample half of the data from in-distribution domains for

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 [17]	TS [14]	MD-TS	MSP [17]	TS [14]	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

212 calibrating models and use the remaining samples for InD ECE evaluation. We use all the samples
 213 from OOD domains for ECE evaluation.

214 **Models and training setup.** We consider multiple network architectures for evaluation, including
 215 ResNet-50 [15], ResNext-50 [46], DenseNet-121 [21], BiT-M-50 [23], Efficientnet-b1 [38], ViT-
 216 Small, and ViT-Base [10]. To evaluate on ImageNet-C, we directly evaluate models that are pre-
 217 trained on ImageNet [9]. For WILDS-RxRx1 and GLDv2, we use the ImageNet pre-trained models
 218 as initialization and apply SGD optimizer to training the models on training datasets.

219 **Evaluation metrics.** We use the Expected Calibration Error (ECE) as the main evaluation metric.
 220 We set the bin size as 100 for ImageNet-C, and set bin size as 20 for WILDS-RxRx1 and GLDv2.
 221 We evaluate ECE on both InD domains and OOD domains. Specifically, we evaluate the ECE of
 222 each InD/OOD domain. Meanwhile, we also evaluate the ECE of the pooled InD/OOD domains, i.e.,
 223 the ECE evaluated on all samples from InD/OOD domains. We use unseen samples from the InD
 224 domain to measure the per-domain ECE. We also measure the averaged per-domain ECE results (i.e.,
 225 per-domain ECE averaged across domains).

226 4.1 Main results

227 We summarize the ECE results of different methods on three datasets in Table 1 and Figure 3. We
 228 use TS to denote temperature scaling [14], and use MSP to denote applying the maximum softmax
 229 probability [17] of the model output (i.e., without calibration). In Table 1, we use the ImageNet
 230 validation dataset and ImageNet-C datasets with severity level $s \in \{1, 5\}$ as the InD domains and use
 231 the remaining datasets as OOD domains. We present the averaged per-domain ECE results in Table 1,
 232 and visualize the ECE of each domain in Figure 3. As shown in Table 1 and Figure 3(a)-3(c), we
 233 find that our proposed approach achieves much better InD per-domain calibration compared with
 234 baselines. Also, TS does not significantly improve over MSP on ImageNet-C InD domains in Table 1,
 235 but our proposed method largely improve the ECE compared with MSP and TS. For instance, the
 236 ECE results of MSP and TS on Efficientnet-b1 are 6.93 and 6.99, and our method achieves 3.84.
 237 Intuitively, when there are a diverse set of domains in the calibration dataset, a single temperature
 238 cannot provide well-calibrated confidences. In contrast, our proposed method is able to produce
 239 much better InD confidence estimates by leveraging the multi-domain structure of the data.

240 Next we study the performance of different methods on out-of-distribution domains. From Table 1,
 241 we find that MD-TS achieves the best performance on OOD domains across all the settings. On
 242 ImageNet-C with BiT-M-R50, MD-TS improves the ECE from 6.54 (MSP) to 4.05, while the
 243 performance of TS is similar to MSP. Moreover, MD-TS significantly outperforms MSP and TS on
 244 WILDS-RxRx1, where MD-TS improves over TS by around 5.00 measured in ECE. Figure 3(d)-3(f)
 245 display the per-domain ECE performance on out-of-distribution domains. MD-TS improves over TS
 246 on more than half of the domains in all three datasets. For the remaining domains, MD-TS performs
 247 slightly worse than TS. Furthermore, on those domains that TS performs poorly (ECE > 8), MD-TS
 248 largely improves over TS by large margins. Further comparisons in Appendix B.7 show that these

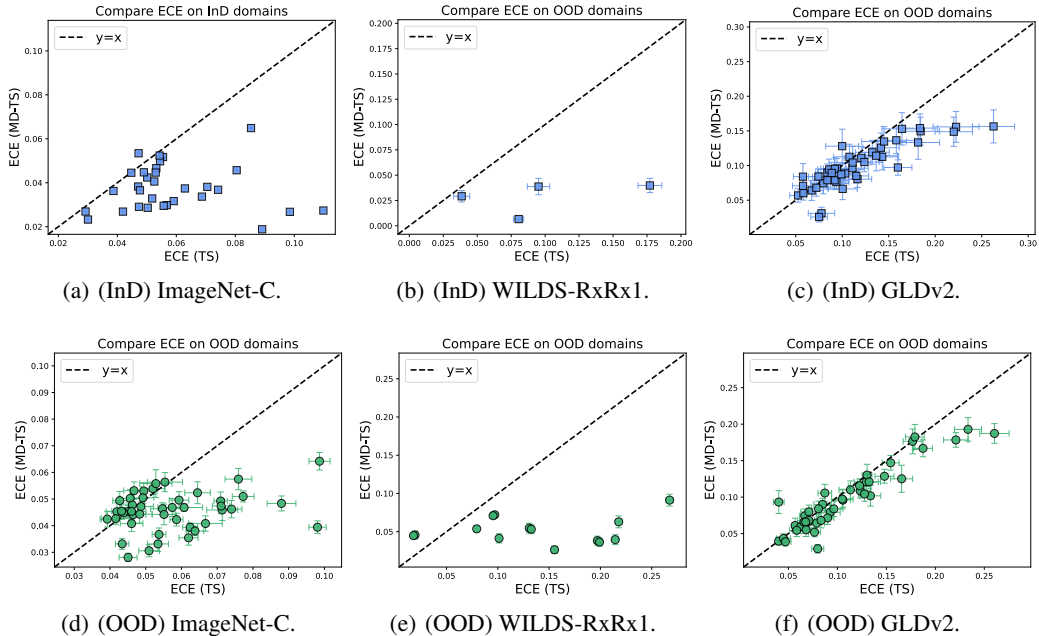


Figure 3: Per-domain ECE of MD-TS and TS on both in-distribution domains and out-of-distribution domains. Each plot is shown with ECE of TS (X -axis) and ECE of MD-TS (Y -axis). Top: per-domain ECE evaluated on InD domains. Bottom: per-domain ECE evaluated on OOD domains. Lower ECE is better.

249 [improvements continue to hold when relative to two other calibration techniques: MC dropout \[11\]](#)
 250 [and deep ensembles \[26\].](#)

251 4.2 Predicting generalization

252 Suppose a model can produce calibrated confidences on unseen samples, in which case we could
 253 leverage the calibrated confidence to predict the model performance. Specifically, based on the
 254 definition of ECE in Eq. (1), when the model is well-calibrated, the average of the calibrated
 255 confidence is close to the average accuracy, i.e., $\text{Conf}(\mathcal{D}) \approx \text{Acc}(\mathcal{D})$.² Meanwhile, predicting model
 256 performance accurately is an essential ingredient in developing reliable machine learning systems,
 257 especially under distributional shifts [13]. As shown in Table 1, we find that our proposed method
 258 produces well-calibrated confidence values on both InD and OOD domains. We now measure its
 259 performance on predicting model performance and compare with existing methods. We measure
 260 the performance using mean absolute error (MAE), $\text{MAE} = (1/K) \cdot \sum_{k=1}^K |\text{Conf}(\mathcal{D}_k) - \text{Acc}(\mathcal{D}_k)|$
 261 where S_k is the dataset from the k -th domain.

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

269 4.3 MD-TS ablations

270 To learn a calibration model that performs well per-domain, we apply linear regression on feature
 271 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-
 272 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} .

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 [17]	TS [14]	MD-TS	MSP [17]	TS [14]	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

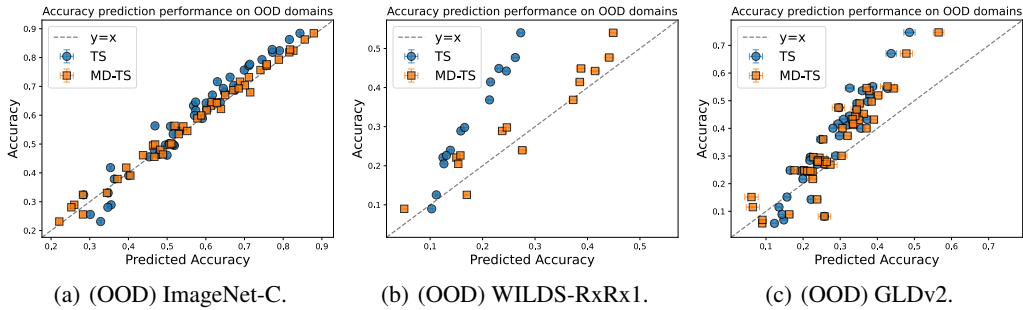


Figure 4: Predicting accuracy performance of MD-TS and TS on both out-of-distribution domains. Each plot is shown with predicted accuracy (X -axis) and accuracy (Y -axis). Each points corresponds to one domain. The network architecture is ResNet-50 for three datasets. Point closer to the $Y = X$ dashed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

273 representations to temperatures in a regression framework. Specifically, beside the ordinary least
 274 squares (OLS) used in Algorithm 1, we consider ridge regression (Ridge), robust regression with
 275 Huber loss (Huber), kernel ridge regression (KRR), and K -nearest neighbors regression (KNN). The
 276 implementations are mainly based on `scikit-learn` [34]. We use grid search (on InD domains) to
 277 select hyperparameters for Ridge, Huber, KRR, and KNN.

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

284 **Efficiency.** We study the efficiency of our proposed method by measuring running time on
 285 three datasets (seconds). We consider the ResNet50 for all datasets. By using the standard
 286 `sklearn.linear_model`, it takes 39.8s/3.2s/4.1s on ImageNet-C/WILDS-RxRx1/GLDv2 for solving
 287 the linear regression problem of MD-TS, and the overall running time of MD-TS is 49.7s/4.6s/6.6s
 288 on ImageNet-C/WILDS-RxRx1/GLDv2. standard TS takes 7.8s/1.2s/3.6s on ImageNet-C/WILDS-
 289 RxRx1/GLDv2. We summarize the comparison results in Table 9, Appendix B.9.

290 5 Theoretical analysis

291 In this section, we provide theoretical analysis to support our understanding of our proposed algorithm
 292 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
 293 calibration map for the base model f on the k -th domain; this map *minimizes* the expected calibration
 294 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
 295 hypothesis class \mathcal{H} . Next, given the fixed base model f , we aim to learn $\hat{h}(\cdot) = h(\cdot; f, \hat{\beta})$ such
 296 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

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

297 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
 298 generalizing to new domains: suppose there is an unseen OOD domain $\tilde{\mathcal{P}}$ and its marginal feature
 299 distribution is different from existing domains, i.e., $\tilde{\mathcal{P}}_X \neq \mathcal{P}_{k,X}$ for $k \in [K]$.

300 Our goal is to understand the conditions under which \hat{h} can have similar calibration on OOD domains
 301 as the InD domains. For example, if the OOD domain is similar to the mixture distribution of InD
 302 domains, we would expect \hat{h} performs similarly on InD and OOD domains. To quantify the distance
 303 between two distributions, we first introduce the \mathcal{H} -divergence [5] to measure the distance between
 304 two distributions:

305 **Definition 5.1** (\mathcal{H} -divergence). *Given an input space \mathcal{X} and two probability distributions \mathcal{P}_X and*
 306 *\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*
 307 *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*
 308 *\mathcal{P}'_X is defined as*

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

309 The \mathcal{H} -divergence reduces to the standard total variation (TV) distance when \mathcal{H} contains all mea-
 310 surable functions on \mathcal{X} , which implies that the \mathcal{H} -divergence is upper bounded by the TV-distance,
 311 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
 312 dimension or pseudo-dimension, the \mathcal{H} -divergence can be estimated using finite samples from \mathcal{P}_X
 313 and \mathcal{P}'_X [5]. Next, we define the mixture distribution of the K in-distribution domains $\mathcal{P}_{K,X}^\alpha$ on input
 314 space \mathcal{X} as follows:

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

315 Given multiple domains $\{\mathcal{P}_1, \dots, \mathcal{P}_K\}$, we can optimize the combination parameters α such that
 316 $\mathcal{P}_{K,X}^\alpha$ minimizes the \mathcal{H} -divergence between $\mathcal{P}_{K,X}^\alpha$ and $\tilde{\mathcal{P}}_X$. More specifically, we define $\hat{\alpha}$ as

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

317 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 $\tilde{\mathcal{H}} := \{\operatorname{sign}(|h(x) - h'(x)| -$
 318 $t) \mid h, h' \in \mathcal{H}, 0 \leq t \leq 1\}$. We now give an upper bound on the risk on the unseen OOD domain. This
 319 result follows very closely those of Blitzer et al. [6], Zhao et al. [49], instantiated in our calibration
 320 setup. Details can be found in Appendix C.

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

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

325 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
 326 OOD domain, $\text{Pdim}(\mathcal{H})$ is the pseudo-dimension of the hypothesis class \mathcal{H} , and $\hat{\varepsilon}(\hat{h}, \mathcal{D}_{k,X})$ is the
 327 empirical risk of the hypothesis \hat{h} on $\mathcal{D}_{k,X}$.

328 **Remark 5.3.** *As shown in Theorem 5.2, even if the OOD domain is very different from the in-*
 329 *distribution domains, the Eq. (4) still implies that we could decrease the risk upper bound on the*
 330 *OOD domain if we perform multi-domain calibration. More specifically, if we could achieve good*
 331 *calibration performance on each individual domain by using multi-domain calibration (which that*
 332 *the first term in the RHS of Eq. (4) is small), then the term $\frac{1}{2}d_{\mathcal{H}}(\mathcal{P}_{K,X}^{\hat{\alpha}}, \tilde{\mathcal{P}}_X) + \lambda(\mathcal{P}_{K,X}^{\hat{\alpha}}, \tilde{\mathcal{P}}_X)$ is*
 333 *always smaller or equal to $\frac{1}{2}d_{\mathcal{H}}(\mathcal{P}', \tilde{\mathcal{P}}_X) + \lambda(\mathcal{P}', \tilde{\mathcal{P}}_X)$, where \mathcal{P}' is the pooled distribution or any*
 334 *individual domain distribution.*

335 **Remark 5.4.** *As suggested by Eq. (4) of Theorem 5.2, larger risks on in-distribution domains will*
 336 *lead to a larger upper bound for the risk evaluated on the OOD domain. On the other hand, as shown*
 337 *in Figure 1, a universal temperature is not sufficient to achieve good calibration performance on each*
 338 *individual in-distribution domain. Therefore, even in the mixture of in-distribution domain setting, a*
 339 *universal temperature is suboptimal and applying multi-domain temperature scaling could be better*
 340 *than using a universal temperature.*

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

348 6 Discussion

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

355 7 Societal Impact

356 In this paper, we aim to improve the trustworthiness of machine learning systems by first, explicitly
 357 accounting for uncertainty, and second, do this in a way that is robust to distribution shifts. As
 358 uncertainty quantification is an increasingly important component of real-world machine learning
 359 systems, including health care and autonomous driving, we believe our work could potentially benefit
 360 a wide range of societal activities. Moreover, our method explicitly takes into account subgroups
 361 of the data, trying to achieve good performance across all data. It is known that this is an important
 362 aspect of performance when deploying models in high-consequence settings [2]. We hope our work
 363 could offer a new perspective on uncertainty quantification under distributional shifts. We do not
 364 anticipate the negative social impact of this work.

References

- [1] Anastasios N. Angelopoulos, Stephen Bates, Emmanuel J. Candès, Michael I. Jordan, and Lihua Lei. Learn then test: Calibrating predictive algorithms to achieve risk control. *arXiv preprint*, 2021. arXiv:2110.01052.
- [2] Solon Barocas, Moritz Hardt, and Arvind Narayanan. *Fairness and Machine Learning*. fairml-book.org, 2019. <http://www.fairmlbook.org>.
- [3] Stephen Bates, Anastasios Angelopoulos, Lihua Lei, Jitendra Malik, and Michael I. Jordan. Distribution-free, risk-controlling prediction sets. *Journal of the ACM*, 68(6), September 2021. doi: 10.1145/3478535.
- [4] Sara Beery, Grant Van Horn, and Pietro Perona. Recognition in terra incognita. In *Proceedings of the European conference on computer vision (ECCV)*, pages 456–473, 2018.
- [5] Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman Vaughan. A theory of learning from different domains. *Machine learning*, 79(1): 151–175, 2010.
- [6] John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman. Learning bounds for domain adaptation. *Advances in neural information processing systems*, 20, 2007.
- [7] Maxime Cauchois, Suyash Gupta, Alnur Ali, and John C Duchi. Robust validation: Confident predictions even when distributions shift. *arXiv preprint arXiv:2008.04267*, 2020.
- [8] Jeffrey Dean, Greg Corrado, Rajat Monga, Kai Chen, Matthieu Devin, Mark Mao, Marc’auelio Ranzato, Andrew Senior, Paul Tucker, Ke Yang, et al. Large scale distributed deep networks. *Advances in neural information processing systems*, 25, 2012.
- [9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [10] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2020.
- [11] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059. PMLR, 2016.
- [12] Yarin Gal, Jiri Hron, and Alex Kendall. Concrete dropout. *Advances in neural information processing systems*, 30, 2017.
- [13] Devin Guillory, Vaishaal Shankar, Sayna Ebrahimi, Trevor Darrell, and Ludwig Schmidt. Predicting with confidence on unseen distributions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1134–1144, 2021.
- [14] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *International Conference on Machine Learning*, pages 1321–1330. PMLR, 2017.
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [16] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *arXiv preprint arXiv:1903.12261*, 2019.
- [17] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *arXiv preprint arXiv:1610.02136*, 2016.

- 410 [18] Dan Hendrycks, Mantas Mazeika, Saurav Kadavath, and Dawn Song. Using self-supervised
411 learning can improve model robustness and uncertainty. *Advances in Neural Information*
412 *Processing Systems*, 32, 2019.
- 413 [19] Dan Hendrycks, Norman Mu, Ekin D Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshmi-
414 narayanan. Augmix: A simple data processing method to improve robustness and uncertainty.
415 *arXiv preprint arXiv:1912.02781*, 2019.
- 416 [20] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo,
417 Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, Dawn Song, Jacob Steinhardt, and Justin
418 Gilmer. The many faces of robustness: A critical analysis of out-of-distribution generalization.
419 *ICCV*, 2021.
- 420 [21] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected
421 convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern*
422 *recognition*, pages 4700–4708, 2017.
- 423 [22] Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay
424 Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, et al.
425 Wilds: A benchmark of in-the-wild distribution shifts. In *International Conference on Machine*
426 *Learning*, pages 5637–5664. PMLR, 2021.
- 427 [23] Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain
428 Gelly, and Neil Houlsby. Big transfer (bit): General visual representation learning. In *European*
429 *conference on computer vision*, pages 491–507. Springer, 2020.
- 430 [24] Lingkai Kong, Jimeng Sun, and Chao Zhang. Sde-net: Equipping deep neural networks with
431 uncertainty estimates. In *International Conference on Machine Learning*, pages 5405–5415.
432 PMLR, 2020.
- 433 [25] Aviral Kumar and Sunita Sarawagi. Calibration of encoder decoder models for neural machine
434 translation. *arXiv preprint arXiv:1903.00802*, 2019.
- 435 [26] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable
436 predictive uncertainty estimation using deep ensembles. *Advances in neural information*
437 *processing systems*, 30, 2017.
- 438 [27] Donghwan Lee, Xinmeng Huang, Hamed Hassani, and Edgar Dobriban. T-cal: An optimal test
439 for the calibration of predictive models. *arXiv preprint arXiv:2203.01850*, 2022.
- 440 [28] Matthias Minderer, Josip Djolonga, Rob Romijnders, Frances Hubis, Xiaohua Zhai, Neil
441 Houlsby, Dustin Tran, and Mario Lucic. Revisiting the calibration of modern neural networks.
442 *Advances in Neural Information Processing Systems*, 34, 2021.
- 443 [29] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of machine learning*.
444 2018.
- 445 [30] Mahdi Pakdaman Naeini, Gregory Cooper, and Milos Hauskrecht. Obtaining well calibrated
446 probabilities using bayesian binning. In *Twenty-Ninth AAAI Conference on Artificial Intelligence*,
447 2015.
- 448 [31] Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, David Sculley, Sebastian Nowozin, Joshua
449 Dillon, Balaji Lakshminarayanan, and Jasper Snoek. Can you trust your model’s uncertainty?
450 evaluating predictive uncertainty under dataset shift. *Advances in neural information processing*
451 *systems*, 32, 2019.
- 452 [32] Sangdon Park, Osbert Bastani, James Weimer, and Insup Lee. Calibrated prediction with
453 covariate shift via unsupervised domain adaptation. In *International Conference on Artificial*
454 *Intelligence and Statistics*, pages 3219–3229. PMLR, 2020.
- 455 [33] Sangdon Park, Edgar Dobriban, Insup Lee, and Osbert Bastani. Pac prediction sets under
456 covariate shift. *arXiv preprint arXiv:2106.09848*, 2021.

- 457 [34] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel,
458 P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher,
459 M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine*
460 *Learning Research*, 12:2825–2830, 2011.
- 461 [35] John Platt et al. Probabilistic outputs for support vector machines and comparisons to regularized
462 likelihood methods. *Advances in large margin classifiers*, 10(3):61–74, 1999.
- 463 [36] Aleksandr Podkopaev and Aaditya Ramdas. Distribution-free uncertainty quantification for
464 classification under label shift. In *Uncertainty in Artificial Intelligence*, pages 844–853. PMLR,
465 2021.
- 466 [37] Joaquin Quiñero-Candela, Masashi Sugiyama, Anton Schwaighofer, and Neil D Lawrence.
467 *Dataset shift in machine learning*. Mit Press, 2008.
- 468 [38] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural
469 networks. In *International conference on machine learning*, pages 6105–6114. PMLR, 2019.
- 470 [39] Sunil Thulasidasan, Gopinath Chennupati, Jeff A Bilmes, Tanmoy Bhattacharya, and Sarah
471 Michalak. On mixup training: Improved calibration and predictive uncertainty for deep neural
472 networks. *Advances in Neural Information Processing Systems*, 32, 2019.
- 473 [40] Ryan J Tibshirani, Rina Foygel Barber, Emmanuel Candes, and Aaditya Ramdas. Conformal
474 prediction under covariate shift. *Advances in neural information processing systems*, 32, 2019.
- 475 [41] Vladimir Vovk, Alexander Gammerman, and Glenn Shafer. *Algorithmic learning in a random*
476 *world*. Springer Science & Business Media, 2005.
- 477 [42] Yoav Wald, Amir Feder, Daniel Greenfeld, and Uri Shalit. On calibration and out-of-domain
478 generalization. *Advances in Neural Information Processing Systems*, 34, 2021.
- 479 [43] Ximei Wang, Mingsheng Long, Jianmin Wang, and Michael Jordan. Transferable calibration
480 with lower bias and variance in domain adaptation. *Advances in Neural Information Processing*
481 *Systems*, 33:19212–19223, 2020.
- 482 [44] Tobias Weyand, Andre Araujo, Bingyi Cao, and Jack Sim. Google landmarks dataset v2-a large-
483 scale benchmark for instance-level recognition and retrieval. In *Proceedings of the IEEE/CVF*
484 *conference on computer vision and pattern recognition*, pages 2575–2584, 2020.
- 485 [45] Cihang Xie, Mingxing Tan, Boqing Gong, Jiang Wang, Alan L Yuille, and Quoc V Le. Adver-
486 sarial examples improve image recognition. In *Proceedings of the IEEE/CVF Conference on*
487 *Computer Vision and Pattern Recognition*, pages 819–828, 2020.
- 488 [46] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual
489 transformations for deep neural networks. In *Proceedings of the IEEE conference on computer*
490 *vision and pattern recognition*, pages 1492–1500, 2017.
- 491 [47] Bianca Zadrozny and Charles Elkan. Obtaining calibrated probability estimates from decision
492 trees and naive bayesian classifiers. In *Icml*, volume 1, pages 609–616. Citeseer, 2001.
- 493 [48] Bianca Zadrozny and Charles Elkan. Transforming classifier scores into accurate multiclass
494 probability estimates. In *Proceedings of the eighth ACM SIGKDD international conference on*
495 *Knowledge discovery and data mining*, pages 694–699, 2002.
- 496 [49] Han Zhao, Shanghang Zhang, Guanhang Wu, José MF Moura, Joao P Costeira, and Geoffrey J
497 Gordon. Adversarial multiple source domain adaptation. *Advances in neural information*
498 *processing systems*, 31, 2018.

499 **Checklist**

- 500 1. For all authors...
- 501 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
502 contributions and scope? [Yes]
- 503 (b) Did you describe the limitations of your work? [Yes]
- 504 (c) Did you discuss any potential negative societal impacts of your work? [No]
- 505 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
506 them? [Yes]
- 507 2. If you are including theoretical results...
- 508 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See Theo-
509 rem 5.2.
- 510 (b) Did you include complete proofs of all theoretical results? [Yes] See Appendix C.
- 511 3. If you ran experiments...
- 512 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
513 mental results (either in the supplemental material or as a URL)? [Yes] See supplemen-
514 tal material.
- 515 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
516 were chosen)? [Yes] See Section 4 and Appendix A.
- 517 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
518 ments multiple times)? [Yes] See Section 4.
- 519 (d) Did you include the total amount of compute and the type of resources used (e.g., type
520 of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix A.
- 521 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 522 (a) If your work uses existing assets, did you cite the creators? [Yes] See Section 4.
- 523 (b) Did you mention the license of the assets? [N/A]
- 524 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
525
- 526 (d) Did you discuss whether and how consent was obtained from people whose data you’re
527 using/curating? [N/A]
- 528 (e) Did you discuss whether the data you are using/curating contains personally identifiable
529 information or offensive content? [N/A]
- 530 5. If you used crowdsourcing or conducted research with human subjects...
- 531 (a) Did you include the full text of instructions given to participants and screenshots, if
532 applicable? [N/A]
- 533 (b) Did you describe any potential participant risks, with links to Institutional Review
534 Board (IRB) approvals, if applicable? [N/A]
- 535 (c) Did you include the estimated hourly wage paid to participants and the total amount
536 spent on participant compensation? [N/A]