

Adaptive Calibrator Ensemble: Navigating Test Set Difficulty in Out-of-Distribution Scenarios

Yuli Zou^{1*} Weijian Deng^{2*} Liang Zheng²

¹The Hong Kong Polytechnic University ²The Australian National University

yuli1027.zou@connect.polyu.hk firstnama.lastname@anu.edu.au

Abstract

Model calibration usually requires optimizing some parameters (e.g., temperature) w.r.t an objective function like negative log-likelihood. This work uncovers a significant but overlooked aspect that the objective function is influenced by calibration set difficulty: the ratio of misclassified to correctly classified samples¹. If a test set has a drastically different difficulty level from the calibration set, a phenomenon out-of-distribution (OOD) data often exhibit: the optimal calibration parameters of the two datasets would be different, rendering an optimal calibrator on the calibration set suboptimal on the OOD test set and thus degraded calibration performance. With this knowledge, we propose a simple and effective method named adaptive calibrator ensemble (ACE) to calibrate OOD datasets whose difficulty is usually higher than the calibration set. Specifically, two calibration functions are trained, one for in-distribution data (low difficulty), and the other for severely OOD data (high difficulty). To achieve desirable calibration on a new OOD dataset, ACE uses an adaptive weighting method that strikes a balance between the two extreme functions. When plugged in, ACE generally improves the performance of a few state-of-the-art calibration schemes on a series of OOD benchmarks. Importantly, such improvement does not come at the cost of the in-distribution calibration performance. Project Website: <https://github.com/insysgroup/Adaptive-Calibrators-Ensemble.git>.

1. Introduction

Model calibration aims to connect the neural network output with uncertainty. A common practice is to find optimal parameters against certain objective functions on a held-out calibration set, to obtain an optimized *calibrator*. In this paper, we focus on post-hoc calibration methods, which

¹Dataset difficulty (w.r.t a classifier) shares the same meaning of classifier accuracy on this dataset. We use “difficulty” to indicate the property of a dataset (i.e., its OOD degree to the classifier), instead of using “accuracy” which describes the performance of the classifier on a dataset.

require training a calibration mapping function to rescale the confidence scores of a trained neural network to make it calibrated [9, 10, 18]. A popular technique is Temperature Scaling [9], which optimizes model temperature by minimizing the negative log-likelihood (NLL) loss.

Post-hoc calibration methods generally work well when calibrating in-distribution test sets. However, oftentimes their calibration performance drops significantly when being tested on an out-of-distribution (OOD) test set [24]. For example, temperature scaling has shown to be ineffective under distribution shift in some scenarios [24]. This problem happens because the test environment (OOD) is different from the training environment due to factors like sample bias and non-stationarity. This paper thus aims to improve post-hoc calibration methods by producing reliable and predictive uncertainty under distribution shifts.

In the community, there exist a few works studying the OOD calibration problem [28, 31, 36]. They typically aim to make amendments to the calibration set to let it approximate the OOD data in certain aspects [28, 31]. Nevertheless, these techniques are typically not adaptive to the test dataset, that is, the calibration set transformation process cannot automatically adjust to the test set. In our experiment, we observe that they improve calibration on some OOD datasets but significantly lead to decreased in-distribution calibration performance. In this regard, while TransCal [36] can perform domain adaptation according to the test domain, it needs to be re-trained for every new test set.

In this paper, our contributions are mainly in two aspects. **First**, we provide a new perspective to understand calibration failure on out-of-distribution datasets. Specifically, we show that **the calibration objective is dependent on the dataset difficulty**. When the calibration set have the same distribution with the test set, it has low difficulty, and thus the calibrator learned on the calibration set would be effective on the test set [9, 10, 18]. Yet, out-of-distribution test sets usually exhibit a different (usually higher) difficulty level compared with the calibration set because of the distribution gap. Under this case, the optimal calibration functions are different for calibration and OOD datasets. That is, a calibra-

tor that optimized on the calibration set would not be optimal on OOD data and thus achieve poor calibration performance.

Second, to achieve robust calibration under distribution shifts, we propose a simple but effective method named adaptive calibrator ensemble (ACE). It adaptively integrates two predefined calibrators: 1) one trained on an easy in-distribution dataset, and 2) the other trained on a severely OOD data set with high difficulty. By estimating how much a new test set deviates from the high-difficulty calibration set, we compute a test adaptive weight to balance the force between the two calibrators. We show that our proposed ACE method improves three existing post-hoc calibration algorithms such as Spline [10] on commonly used OOD benchmarks. Moreover, our method does *not* have compromised calibration performance for in-distribution data.

2. Related Work

Post-hoc calibration calibrates a trained neural network by rescaling confidence scores [1, 9, 10, 16, 18, 21, 23, 26, 30, 34, 39, 40]. For example, as a multi-class extension of Platt scaling, vector scaling and matrix scaling [9] introduce a linear layer to transform the logits vector to calibrate the network outputs. Spline [10] obtain a recalibration function via spline-fitting, which directly maps the classifier outputs to the calibrated probabilities. Dirichlet [18] propose a multi-class calibration method, derived from Dirichlet distributions. Rahimi *et al.* [26] propose a general post-hoc calibration function that can preserve the top- k predictions of any deep network via an intra-order-preserving function. Our work seeks to improve the OOD performance of existing post-hoc calibrators such as vector scaling, temperature scaling, and spline, through an ensemble mechanism.

Out-of-distribution calibration. A few works study calibration under distribution shift [28, 36]. To improve the post-hoc calibration under distribution shift, some studies [28, 31] propose to modify the calibration set to represent a generic distribution shift. Moreover, prediction uncertainty is studied in [17]. Based on the uncertainty, an ‘‘accuracy versus uncertainty’’ calibration loss is proposed to encourage a model to be certain of correctly classified samples and uncertain of inaccurate samples. In comparison, our method is based on *whether samples are correctly or incorrectly classified (i.e., difficulty) rather than uncertainty*. We find difficulty is an important factor for OOD calibration failure. Furthermore, TransCal [36] uses unsupervised domain adaptation to improve temperature scaling. This method has a high computational cost because, 1) it needs an additional domain adaptation training process, and 2) every time it meets a new test set, the domain adaptation model needs to be re-trained. Gong *et al.* [7] study the calibration under domain generalization setting where they develop calibration methods on calibration sets from *multiple* domains. We contribute from a different perspective to the existing liter-

ature. We provide insight into the role of dataset difficulty on the failure of existing algorithms on OOD data. We then propose a simple and effective ensemble strategy to improve post-hoc calibrators in a test set adaptive manner.

3. Methodology

3.1. Preliminaries

Neural network notations. Considering the task of calibrating neural networks for n -way classification, let us define $[n] := \{1, \dots, n\}$, $\mathcal{X} \subseteq \mathbb{R}^d$ be the domain, $\mathcal{Y} = [n]$ be the label space, and Δ_n denote the $n - 1$ dimensional unit simplex. Given a training dataset \mathcal{D}_{tr} of independent and identically distributed (i.i.d.) samples drawn from an unknown distribution π on $\mathcal{X} \times \mathcal{Y}$, we learn a probabilistic predictor $\phi : \mathbb{R}^d \rightarrow \Delta_n$. We assume that ϕ can be expressed as the composition $\phi =: \mathbf{sm} \circ \mathbf{g}$, with $\mathbf{g} : \mathbb{R}^d \rightarrow \mathbb{R}^n$ being a non-probabilistic n -way classifier and $\mathbf{sm} : \mathbb{R}^n \rightarrow \Delta_n$ being the softmax operator $\mathbf{sm}_i(\mathbf{z}) = \frac{\exp(\mathbf{z}_i)}{\sum_{j=1}^n \exp(\mathbf{z}_j)}$, for $i \in \mathcal{Y}$, where the subscript i denotes the i -th element of a vector. We say $\mathbf{g}(\mathbf{x})$ is the *logits* of \mathbf{x} with respect to ϕ .

Definition of a calibrated network. When queried at $(\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y}$ sampled from an unknown distribution π , the probabilistic predictor ϕ returns $\hat{y} =: \arg \max_i \phi_i(\mathbf{x})$ as the predicted label and $\hat{p} =: \max_i \phi_i(\mathbf{x})$ as the associated confidence score. We say ϕ is *perfectly calibrated* with respect to π , if \hat{p} is expected to represent the true probability of correctness. Formally, a perfectly calibrated model satisfies $\mathbb{P}(\hat{y} = y | \hat{p} = p) = p$ for any $p \in [0, 1]$. In practice, we commonly use the Expected Calibration Error (ECE) [9] as the calibration performance metric. It first groups all samples into M equally interval bins $\{B_m\}_{m=1}^M$ with respect to their confidence scores, and then calculates the expected difference between the accuracy and average confidence: $\text{ECE} = \sum_{m=1}^M \frac{|B_m|}{n} |\text{acc}(B_m) - \text{avgConf}(B_m)|$, where n denotes the number of samples.

Post-hoc calibration learns a post-hoc calibration function $\mathbf{f} : \mathcal{R}^n \rightarrow \mathcal{R}^n$ such that the new probabilistic predictor $\phi_c := \mathbf{sm} \circ \mathbf{f} \circ \mathbf{g}$ is better calibrated *and* tries to keep a similar (or same) accuracy of the original network ϕ .

3.2. Post-hoc Calibration Function Is Influenced by Calibration Set Difficulty

Post-hoc calibration loss function. Assume we have a held-out calibration dataset $\mathcal{D}_c = \{(\mathbf{x}^i, y^i)\}_{i=1}^N$ with i.i.d samples from the unknown distribution π on $\mathcal{X} \times \mathcal{Y}$ and a calibration function \mathbf{f} parameterized by some vector θ . The empirical calibration loss is generally defined as,

$$\frac{1}{N} \sum_{i=1}^N \ell(y^i, \mathbf{f}(\mathbf{z}^i)) + \frac{\lambda}{2} \|\theta\|^2, \quad (1)$$

where $\mathbf{z}^i = \mathbf{g}(\mathbf{x}^i)$, $\ell : \mathcal{Y} \times \mathcal{R}^n \rightarrow \mathcal{R}$ is a cost function,

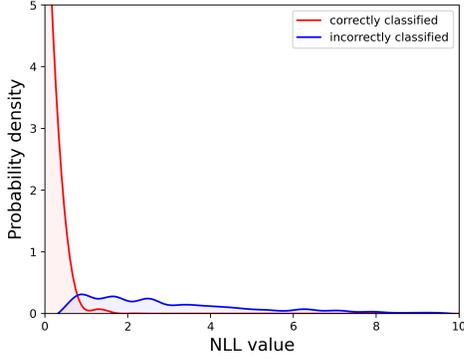


Figure 1. **NLL values of correctly and incorrectly classified samples.** We use ResNet-152 on the in-distribution ImageNet calibration set (as described in Section 4.1) and plot NLL probability density of the two types of samples. We clearly observe that correctly classified samples generally have a much lower NLL value.

and $\lambda \geq 0$ is the regularization weight. $\ell(\cdot, \cdot)$ is the network classification loss. Following existing literature, we employ the commonly used negative log-likelihood (NLL) loss:

$$\ell(y, \mathbf{f}(\mathbf{z})) = -\log(\mathbf{sm}_y(\mathbf{f}(\mathbf{g}(\mathbf{x}))))), \quad (2)$$

where \mathbf{sm} is softmax operator, and \mathbf{sm}_y is its y -th element.

Plain fact: individual samples matter in the classification loss. Apparently, a major component in the calibration objective (Eq. 1) is the model classification loss (e.g., the commonly used NLL loss, Eq. 2). If a sample is correctly classified, the classification loss will likely return a small value; If a sample is incorrectly classified, there will likely be a high loss value. Therefore, *whether an individual sample is correctly classified or not would lead to quite different classification loss values.*

We conduct an empirical analysis to verify this conclusion. Specifically, we use ResNet-152 trained on ImageNet [4]. The NLL values of these samples are computed on the calibration set (described in Section 4.1), and summararily drawn in Fig. 1. It is clearly shown that the NLL values of correctly classified samples are close to 0 while those of incorrectly classified samples are significantly greater.

Collectively, calibration set difficulty influences calibration optimization. To illustrate this point, we use NLL as an example, which is a commonly used classification loss. Given that the two types of samples have different NLL values, we decompose the NLL loss into two parts:

$$\ell_T(y, \mathbf{f}(\mathbf{z})) = -\frac{1}{N_T} \sum_i^{N_T} \log(\mathbf{sm}_{y^i}(\mathbf{f}(\mathbf{g}(\mathbf{x}^i))))), \quad (3)$$

where $\arg \max \mathbf{sm}(\mathbf{g}(\mathbf{x}^i)) = y^i$, and,

$$\ell_F(y, \mathbf{f}(\mathbf{z})) = -\frac{1}{N_F} \sum_i^{N_F} \log(\mathbf{sm}_{y^i}(\mathbf{f}(\mathbf{g}(\mathbf{x}^i))))), \quad (4)$$

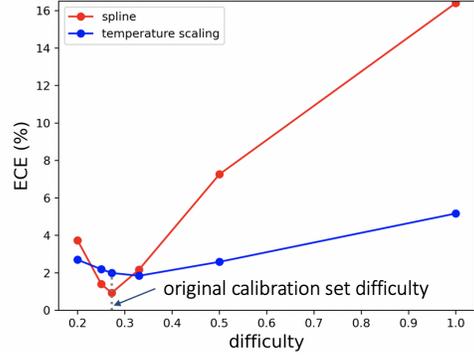


Figure 2. **The impact of calibration set difficulty ($\frac{N_F}{N_T}$) on calibration performance (ECE).** We manually select images from the in-distribution calibration set to create new calibration sets of various difficulty levels. The difficulty of the original calibration set is 0.2723. We use the ResNet-152 model and an in-distribution test set from ImageNet. We evaluate two methods (temperature scaling, and Spline), and at the same time, mark the difficulty of the original in-distribution calibration set (gray vertical dotted line). We find the calibration sets having similar difficulty to the original will lead to good calibration performance and vice versa.

where $\arg \max \mathbf{sm}(\mathbf{g}(\mathbf{x}^i)) \neq y^i$. In Eq. 3 and Eq. 4, N_T and N_F note the numbers of correctly and incorrectly classified samples, respectively. By adjusting $\frac{N_F}{N_T}$, the overall NLL value changes, which will affect the optimized calibration parameters θ (a.k.a. the calibration function).

Formally, we define the *difficulty* of a dataset as $\frac{N_F}{N_T}$. Note that, the difficulty of a dataset (with respect to a classifier) shares the same meaning as classifier accuracy on this dataset. The above analysis indicates that *optimized calibration parameters are affected by the difficulty of the calibration set*: **1)** θ trained on a more difficult calibration set tends to have a larger classification loss values (Eq. 2) and thus a larger calibration loss (Eq. 1). **2)** θ trained on an easier calibration set likely corresponds to a smaller classification loss (Eq. 2) and thus a lower calibration loss (Eq. 1).

We empirically verify the above conclusion in Fig. 2, where we create calibration sets with various levels of difficulty ($\frac{N_F}{N_T}$) and mark the difficulty level of the original calibration set. It indicates that calibration set difficulty indeed influences ECE of two calibration methods: temperature scaling (NLL) [9], Spline (KS-error) [10]. Moreover, when testing the original in-distribution data, if the difficulty of the created calibration dataset is similar, the two calibration methods generally have good calibration performance. However, calibration performance is poorer when the *difficulty* of created calibration dataset is very different from that of the original calibration dataset.

The above analysis mostly uses the NLL loss as an example, but can also apply to some other classification loss functions (e.g., the KS-error used in Spline is verified in

Fig. 2, the cross-entropy loss and the focal loss). These loss functions are usually influenced by individual samples, and thus collectively the dataset difficulty would eventually impact the calibration performance.

3.3. Calibration Set Difficulty Influences Out-of-distribution Calibration

Having analyzed that calibration set difficulty influences the calibration performance on in-distribution test sets, we provide a tentative explanation of *why calibrators trained on in-distribution data fail on OOD test sets*. Essentially, an OOD test set usually has a different difficulty level from the in-distribution calibration set. In fact, the OOD difficulty level is usually higher, *i.e.*, there is a higher percentage of incorrectly classified samples, because of the domain gap problem [2, 5, 22]. Therefore, the calibration mapping function \mathbf{f} that an OOD test set needs is different from the in-distribution calibration set. When training a calibrator on in-distribution data, its performance would thus be suboptimal on the OOD test data.

Moreover, our reasoning also helps understand why some existing OOD calibration methods have compromised calibration performance on in-distribution test sets. Specifically, these methods (*e.g.*, Perturbation [31]) obtain their mapping functions on some modified in-distribution calibration set (*e.g.*, adding Gaussian noise), which to some extent mimics the OOD test set. However, this modification operation is not adaptive, that is, they do not change *w.r.t* the test set. When the test set changes to an in-distribution one, its optimal calibration parameters would be different from those obtained from the modified calibration set. This is possible because of different difficulty levels.

3.4. Adaptive Calibrator Ensemble

Overview. To achieve desirable calibration under distribution shifts, we propose a simple and effective method called Adaptive Calibrator Ensemble (ACE). Using an in-distribution calibration set as input, ACE outputs an OOD calibrator as if having been trained on a calibration set with a proper difficulty level. To do so, we first seek two calibration sets with extreme difficulty levels: an in-distribution difficulty level (easy) and a high difficulty level (hard). We then use an adaptive weighting scheme to fuse the output of calibrators trained on the two extreme calibration sets.

Finding two datasets with extreme difficulty levels. Straightforwardly, we secure the “easy” one as the in-distribution calibration set itself \mathcal{D}_o . To obtain the “hard” calibration set \mathcal{D}_h , we perform sampling on \mathcal{D}_o aiming to increase the difficulty. Specifically, we apply the classifier on the in-distribution calibration set to find correctly classified samples, incorrectly classified samples, and their numbers N_T^o and N_F^o (N_T^o is usually greater than N_F^o). To create \mathcal{D}_h , we calculate its N_T and N_F as follows,

$N_F^h = N_F^o, N_T^h = N_F^o/d$, where $d \in (0, \infty)$ is a pre-defined difficulty level (hyperparameter). We then randomly sample \mathcal{D}_o to achieve this difficulty level. When d is relatively large², the calibration set contains many more incorrectly classified samples than correctly classified ones, allowing us to have the desired calibration set \mathcal{D}_h , which is considered seriously out-of-distribution and hard.

Training two calibrators on the two extreme datasets.

On *each* of the obtained the easy and the hard calibration sets \mathcal{D}_o and \mathcal{D}_h , we train a calibrator. Let \mathbf{g} denote the deep learning model. For calibration dataset $\mathcal{D}_o = \{(\mathbf{x}^i, y^i)\}_{i=1}^{N_o}$, where N_o means the number of samples of \mathcal{D}_o , we train a calibration function \mathbf{f}_o , and the calibrated logits are denoted as $\mathbf{z}_o^i = \mathbf{f}_o(\mathbf{z}_{\text{ori}}^i)$. Here, $\mathbf{z}_{\text{ori}}^i$ is the original uncalibrated logits for a new test set that is either in-distribution or OOD. Similarly, for calibration set $\mathcal{D}_h = \{(\mathbf{x}^i, y^i)\}_{i=1}^{N_h}$, where N_h means the number of samples of \mathcal{D}_h , we train a calibration function \mathbf{f}_h , the calibrated logits is $\mathbf{z}_h^i = \mathbf{f}_h(\mathbf{z}_{\text{ori}}^i)$.

An adaptive method to ensemble outputs of the two calibrators. Given \mathcal{D}_o and \mathcal{D}_h , we intuitively speculate that the difficulty of a usual out-of-distribution test set would be positioned in between. As such, we propose to compute an adaptive weight α to balance the difficulty of these two outputs produced by calibrators, then the final output \mathbf{z}_{cal} is:

$$\mathbf{z}_{\text{cal}} = \alpha \cdot \mathbf{z}_o + (1 - \alpha) \cdot \mathbf{z}_h. \quad (5)$$

In designing a reasonable weight α , we request it to be test-set-adaptive. **First**, when the distribution of an OOD test set is similar to the original calibration set \mathcal{D}_o , $\alpha \rightarrow 1$, so that the system reduces to in-distribution calibrator \mathbf{z}_o ; **Second**, when a test set is seriously out-of-distribution, $\alpha \rightarrow 0$.

Moreover, average confidence score could serve as an unsupervised indicator of the degree of how out-of-distribution a test set is [8]. So given an unlabeled test set $\mathcal{D}_{\text{test}}$, we can estimate an approximate OOD degree of this test set. Here, we compute the ad-hoc weight α as,

$$\alpha = \frac{\text{avgConf}(\mathcal{D}_{\text{test}})}{\text{avgConf}(\mathcal{D}_o)}, \quad (6)$$

where $\text{avgConf}(\cdot)$ calculates the average confidence score of a dataset. In the experiment, we will evaluate some fixed values of α , which are useful on some occasions but less so on others. Moreover, being fixed implies that it does not work for in-distribution data unless it is fixed to 1.

The ensemble scheme works efficiently. To illustrate how ACE ensembles the two calibrators, here we use Temperature Scaling [9] as an example whose calibration function is $\mathbf{f}(\mathbf{z}) = \mathbf{T} \cdot \mathbf{z}$ where \mathbf{T} is a learnable scalar parameter. Let \mathbf{T}_o and \mathbf{T}_h denote the temperature value which learned

²By default, we set $d = 10$, which means 10 times more incorrectly classified samples than correct ones. Notice that we set $d = 9$ for CIFAR-10-C, which equals the randomly classified result.

on the easy calibration set \mathcal{D}_o and the hard calibration set \mathcal{D}_h , respectively. \mathbf{z}_{ori} is the uncalibrated logits of test set. Referring to Eq. 5, the calibrated logits of test set \mathbf{z}_{cal} is:

$$\begin{aligned} \mathbf{z}_{\text{cal}} &= \alpha \cdot \mathbf{z}_{\text{ori}} \cdot \mathbf{T}_o + (1 - \alpha) \cdot \mathbf{z}_{\text{ori}} \cdot \mathbf{T}_h \\ &= \mathbf{z}_{\text{ori}} \cdot (\alpha \cdot \mathbf{T}_o + (1 - \alpha) \cdot \mathbf{T}_h). \end{aligned} \quad (7)$$

Thus the equivalent value of temperature \mathbf{T}_{cal} which has $\mathbf{z}_{\text{cal}} = \mathbf{z}_{\text{ori}} \cdot \mathbf{T}_{\text{cal}}$ can be computed as:

$$\mathbf{T}_{\text{cal}} = \alpha \cdot \mathbf{T}_o + (1 - \alpha) \cdot \mathbf{T}_h. \quad (8)$$

According to Eq. 8, we show that the output-space ensemble (Eq. 5) is equal to the weight-space ensemble of two calibrators. Additional, weight-space ensemble methods have shown superior performance and robustness gains over single models [11, 19, 25, 37, 41]. Therefore, the outputs produced by our ensemble scheme shows to have better calibration performance than single calibrator produces.

3.5. Discussion

Difficulty is a relative concept. Despite being formulated as $\frac{N_F}{N_T}$, difficulty also depends on the model or classifier. For stronger models, the difficulty level would be lower (even $N_F = 0$) and vice versa. In this paper, we assume fixed models and choose not to put the model as a subscript in the definition of difficulty for simplicity.

Domain gap vs. difficulty. Domain gap is used to describe the distribution difference between domains and certainly exists between an OOD test set and the calibration set. Therefore, a possible way to calibrate OOD data is to find a dataset with similar distribution to the OOD test set, which is essentially reflected in [28, 31]. Note that distribution shift does not equal high difficulty. In fact, high difficulty is one of consequences of distribution shift/OOD data. Our paper points out a new way to craft the domain gap by modifying the difficulty of the calibration set. In fact, domain gap is a complex phenomenon and related to many factors aside from difficulty, so it would be interesting to investigate other factors which can help OOD calibration.

An alternative method. We emphasize the main contribution is to report that calibration set difficulty is influential on OOD calibration performance. The designed method, in comparison, is more from an intuitive perspective. There might be other alternatives. For example, we could use the average confidence of a dataset (we use it in Eq. 6 to calculate α instead) to estimate its difficulty and create a calibration set that has a closer difficulty level to the OOD test dataset. We show this alternative also gives improvement over some baselines. (Please refer to the supplemental material for more details.)

Potential limitation and direction. Our weighting method (Eq. 5) assumes that an OOD test set sits between \mathcal{D}_o and \mathcal{D}_h in terms of difficulty. This assumption should

be valid for most cases in practice because the difficulty of \mathcal{D}_o is very low and that of \mathcal{D}_h is very high (we use $d = 10$ by default, which translates to 9.1% top-1 accuracy). We empirically observe that $d = 10$ is effective, which translates to an accuracy of 9.09%. We believe a dataset with 9.09% accuracy is difficult enough to cover a wide range of test sets. Moreover, distribution shift occurs in a variety of ways [3, 13, 29, 20, 33]. There might exist scenarios (e.g., adversarial attack) where confidence score is less effective in describing the distribution shift. In such cases, our method might not be able to achieve significant improvement over existing algorithms. In fact, it would be interesting to explore other potential ways to characterize distribution discrepancy. Furthermore, in realistic application scenarios, we may have access to calibration datasets from multiple domains [7]. To better use these data, one potential way is to learn a ACE model on each calibration set. Then, we ensemble the results of all learned ACE models for a given unknown test set. We evaluate our proposed ACE method under the domain generalization setting in the supplemental material.

4. Experiment

4.1. Experimental Setup

Neural Networks. We consider both convolutional and non-convolutional networks. Specifically, we use ResNet-152 [12], ViT-Small-Patch32-224 [6] and DeiT-Small-Patch16-224 [32]. The three networks are either trained or fine-tuned on the ImageNet training set [4].

Calibration set and in-distribution test set. Following the protocol in [10], we divide the validation set of ImageNet into two halves: one for the in-distribution test (namely ImageNet-Val), the other for learning calibration methods (namely calibration set \mathcal{D}_o).

Out-of-distribution test sets. In the experiment, we use the following *six real-world* out-of-distribution benchmarks. (i) ImageNet-V2 [27] is a new version of ImageNet test set. It contains three different sets resulting from different sampling strategies: Matched-Frequency (A), Threshold-0.7 (B), and Top-Images (C). Each version has 10,000 images from 1000 classes; (ii) ImageNet-S(ketch) [35] shares the same 1000 classes as ImageNet but all the images are black and white sketches. It contains 50,000 images; (iii) ImageNet-R(ention) [13] contains artificial renditions of ImageNet classes. It has 30,000 images of 200 classes. Following [13], we sub-select the model logits for the 200 classes before computing calibration metrics. (iv) ImageNet-Adv(ersarial) [15] is adversarially selected to be hard for ResNet-50 trained on ImageNet. It has 7,500 samples of 200 classes. As for ImageNet-R, we sub-select the logits for the 200 classes before computing the calibration metric. Moreover, we test on synthetic CIFAR-10-C(orrptions) and ImageNet-C(orrptions) [14].

Table 1. OOD calibration performance of our method (ACE) integrated with three post-hoc methods: vector scaling, temperature scaling (Temp. Scaling), and Spline. ECE (25 bins, %) for top-1 predictions is reported. We use **ResNet-152** on various image classification datasets with *various distribution shifts*. For each column, the lowest number is in **bold** and the second lowest underlined. Our method (ACE) effectively improves the post-hoc methods on 15 out of 18 occasions. $\blacktriangle/\blacktriangledown$ denotes ECE is lower / higher than the post-hoc method when being used alone, with statistical significance (p -value < 0.05) based on the two-sample t-test.

Methods	ImgNet-V2-A	ImgNet-V2-B	ImgNet-V2-C	ImgNet-S	ImgNet-R	ImgNet-Adv
uncalibrated	9.5016	6.2311	4.3117	24.6332	17.8621	50.8544
Vector Scaling	6.8068	4.2184	2.9258	20.3726	14.5037	44.7593
+ ACE	5.6291 $\pm 0.0397 \blacktriangle$	3.7742 $\pm 0.0237 \blacktriangle$	3.1141 $\pm 0.0150 \blacktriangledown$	15.8747 $\pm 0.0252 \blacktriangle$	10.6343 $\pm 0.0356 \blacktriangle$	40.5773 $\pm 0.0491 \blacktriangle$
Temp. Scaling	4.4413	2.7309	1.6831	15.7879	10.4797	42.6302
+ ACE	<u>3.5615</u> $\pm 0.0028 \blacktriangle$	2.5692 $\pm 0.0013 \blacktriangle$	1.7021 $\pm 0.0001 \blacktriangledown$	<u>10.3915</u> $\pm 0.0092 \blacktriangle$	<u>6.7458</u> $\pm 0.0083 \blacktriangle$	<u>38.0651</u> $\pm 0.0114 \blacktriangle$
Spline	4.5321	1.8034	<u>1.3357</u>	19.6392	13.1116	45.3623
+ ACE	2.8201 $\pm 0.0283 \blacktriangle$	<u>2.0235</u> $\pm 0.0154 \blacktriangledown$	1.0550 $\pm 0.0092 \blacktriangle$	6.9264 $\pm 0.0864 \blacktriangle$	6.8533 $\pm 0.0011 \blacktriangle$	31.0926 $\pm 0.0422 \blacktriangle$

Both these two datasets are modified with synthetic perturbations such as blur, pixelation, and compression artifacts at a range of severities. We use 80 different distortions (16 different types with 5 levels of intensity each) which are the same as those in [24].

Post-hoc calibration methods. In the experiment, we validate the effectiveness of ACE by integrating it with the existing calibration methods through which we obtain calibrated logits z (Section 3.4). Specifically, we use vector scaling [9], temperature Scaling [9], and Spline [10] as baseline calibrators, and compare with a recent method Perturbation [31] which is specifically designed for OOD calibration. In addition, we also compare with more existing methods, *i.e.*, Ensemble [19, 24], SVI [38], SVI-AvUC and SVI-AvUTS [17], to show our method competitive.

4.2. Calibration on Out-of-distribution Datasets

ACE improves calibration methods on OOD datasets. We evaluate our method combined with three post-hoc calibrators on six out-of-distribution test sets and compare it with those calibrators used alone. Table 1 shows ECE (using 25 bins) results of ResNet-152. Our ACE is shown to consistently improve the OOD calibration results of the three baseline calibrators in most of the test cases. For example, when calibrating ResNet-152, our method improves temperature scaling by 0.88%, 0.17%, 5.40%, 3.73% and 4.57% decrease in ECE, on ImageNet-V2-A/B, ImageNet-S/R/Adv, respectively. Under the same settings, the ECE of our method is slightly higher (0.019%) than the baseline on the ImageNet-V2-C dataset. We also report other metrics (*e.g.*, Brier Score, KS-Error) in the supplemental material.

ACE works effectively under two other neural networks. To show the effectiveness of our method for different

backbones, we adopt two transformer models (ViT-Small-Patch32-224 and Deit-Small-Patch16-224) as backbones, and experimental settings are the same as those in Table 1. Table 2 indicates that for backbone ViT-Small-Patch32-224 our method reduces ECE of the three baselines on five out of the six OOD test sets. For example, compared with Spline, ECE of our method is 1.82%, 0.28%, 11.13%, 6.16% and 14.53% lower on ImageNet-V2-A/C, ImageNet-S/R/Adv, respectively. On the other hand, Table 2 demonstrates that for the Deit-Small-Patch16-224 backbone, our method is beneficial on all the six OOD test sets. In addition, comparing the *uncalibrated* results of the three backbones, transformer models generally have a lower ECE under OOD test sets. Specifically, *ViT-Small-Patch32-224* is shown to be superior to *Deit-Small-Patch16-224* on four out of six test sets.

Comparison with the existing calibration methods. In Table 3, we compare our method with the state-of-the-art methods, *i.e.*, various variants of AvUC [17] and Ensemble [19], on CIFAR-10-C and ImageNet-C. Following the protocol in [24, 17], we report the results at intensity 5. Our method improves Spline by reducing ECE by 11.18% and 6.70% on CIFAR-10-C and ImageNet-C, respectively. Compared with these methods, our method is competitive on both ImageNet-C and CIFAR-10-C. For example, for CIFAR-10-C, our method achieves 3.21% and 1.10% lower calibration error than SVI-AvUTS and SVI-AvUC, respectively.

4.3. ACE Does Not Compromise ID Calibration

We show ECE results on in-distribution test set (ImageNet-Val) using ResNet-152. We adopt the same three post-hoc calibration baselines and Perturbation [31] for comparison. As shown in Fig. 3, we observe that the post-hoc calibration baselines themselves effectively reduce the ECE

Table 2. OOD calibration performance of our method (ACE) and Perturbation [31] applied on Spline [10]. We use ECE (%) as evaluation metric and report results using three neural networks: ResNet-152 [12] (ResNet), ViT-Small-Patch32-224 [6] (ViT), and Deit-Small-Patch16-224 [32] (Deit). All the other notations and settings are the same with Table 1. Our method improves the calibrator baselines in 16 out of 18 scenarios, while Perturbation has mixed performance.

Models	Methods	ImgNet-V2-A	ImgNet-V2-B	ImgNet-V2-C	ImgNet-S	ImgNet-R	ImgNet-Adv
ResNet	Spline	<u>4.5321</u>	1.8034	<u>1.3357</u>	19.6392	13.1116	45.3623
	+ ACE	2.8201 ▲	<u>2.0235</u> ▼	1.0550 ▲	6.9264 ▲	<u>6.8533</u> ▲	31.0926 ▲
	+ Perturbation	5.4175 ▼	8.2109 ▼	9.3326 ▼	<u>7.9805</u> ▲	2.9171 ▲	<u>32.3677</u> ▲
ViT	Spline	<u>4.7572</u>	1.6859	<u>1.4683</u>	15.9864	12.5494	38.0404
	+ ACE	2.9329 ▲	<u>2.0832</u> ▼	1.1831 ▲	4.8514 ▲	<u>6.3699</u> ▲	<u>23.5147</u> ▲
	+ Perturbation	5.0302 ▼	6.3854 ▼	7.8929 ▼	<u>5.9254</u> ▲	3.7302 ▲	22.5118 ▲
Deit	Spline	5.0289	<u>2.1261</u>	<u>1.3923</u>	20.7714	9.6996	31.3674
	+ ACE	2.4576 ▲	1.6475 ▲	1.3544 ▲	5.6622 ▲	3.6721 ▲	15.7885 ▲
	+ Perturbation	<u>3.3520</u> ▲	2.4547 ▼	2.9461 ▼	<u>15.9003</u> ▲	<u>8.1481</u> ▲	<u>27.9474</u> ▲

Table 3. Method comparison on CIFAR-10-C and ImageNet-C with ResNet-20 and ResNet-50, respectively. Following the protocol in [24], we report mean ECE (10 bins for CIFAR-10-C and 25 bins for ImageNet-C, %) across 16 different types of data shift at intensity 5 with lowest numbers in **bold** and the second lowest underlined. For each row, we compare ACE with the best of the competing ones (*i.e.*, SVI-AvUC) using the two-sample t-test.

Dataset	Uncalibrated	Ensemble [19]	SVI [38]	SVI-AvUTS [17]	SVI-AvUC	Spline [10]	Spline+ACE
CIFAR-10-C	0.1942	0.1611	0.2389	0.1585	<u>0.1374</u>	0.3382	0.1272 ▲
ImageNet-C	0.3151	0.0880	0.1188	0.0800	<u>0.0542</u>	0.1147	0.0477 ▲

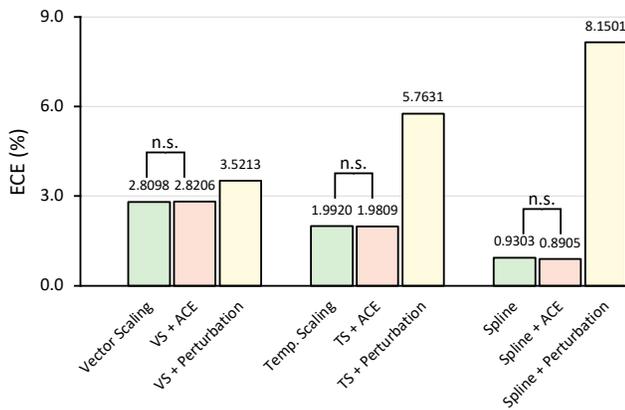


Figure 3. **Evaluation of ACE on the ID test set ImageNet-Val.** We calibrate the ResNet-152 classifier and use ECE (%) for top-1 predictions as evaluation metric. “n.s.” means the difference between results is not statistically significant (p -value $>$ 0.05).

score compared with the uncalibrated system and that Spline generally performs the best. Perturbation is shown to deteriorate the calibration performance for all three baselines. Because Perturbation is not adaptive to different test sets, its effectiveness is not guaranteed when a test set is out of its optimal domain confined by the generated diverse set. In comparison, when our method is integrated with the base-

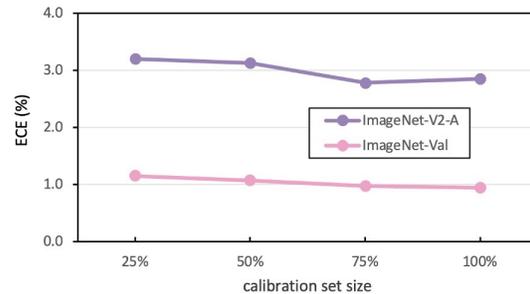


Figure 4. **Effect of the size of the two extreme calibration sets.** Starting from original size (25, 000 and 5, 885 images respectively for \mathcal{D}_o and \mathcal{D}_h), we randomly select a certain percentage of calibration sets. We report ECE of ResNet-152 with Spline on ImageNet-Val and ImageNet-V2-A.

lines, the resulting calibration performance is very close to the baselines when being used alone. This is mainly because of the adaptive weighting scheme (see Section 3.4 for more explanations). Thus our method is *not* compromised on the in-distribution test set.

4.4. Component Analysis of ACE

Impact of the size of the two extreme calibration sets. ACE uses an “easy” calibration set \mathcal{D}_o (the original calibration set) and a “hard” calibration set \mathcal{D}_h . The original \mathcal{D}_o and \mathcal{D}_h have 25, 000 and 5, 885 images, respectively. Here,

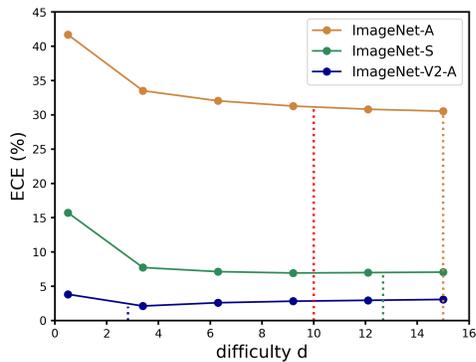


Figure 5. **Impact of hyperparameter d on OOD calibration.** We densely sample values of $d \in (0.5, 15)$ and report ECE (%) of ResNet-152 with Spline on ImageNet-V2-A, ImageNet-S and ImageNet-Adv. We also mark the results using our empirically selected value ($d = 10$) and the optimal values shown by the dotted vertical line with the same color.

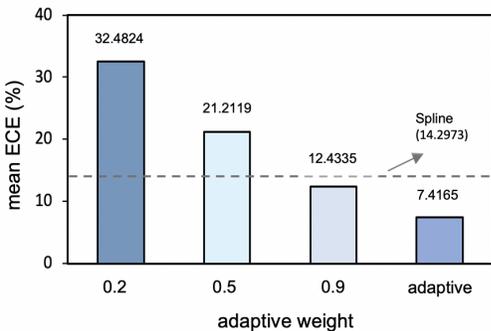


Figure 6. **Comparison of different weighting schemes for ACE.** We report the mean ECE (%) on six OOD datasets (ImageNet-V2-A/B/C, ImageNet-S, ImageNet-Adv, ImageNet-R) and one ID test set (ImageNet-Val). Spline and ResNet-152 is used.

we simultaneously reduce the size of \mathcal{D}_o and \mathcal{D}_h by a certain percentage and report calibration performance (ECE) in Fig. 4. From the results on the in-distribution dataset ImageNet-Val and out-of-distribution dataset ImageNet-V2-A, we observe that our method is relatively stable on both test sets when the size changes. Yet for best results, we recommend using possibly large calibration sets.

Impact of the difficulty of \mathcal{D}_h . To analyze the impact of hyperparameter d (Section 3.4), we create multiple \mathcal{D}_h with various values of d . Results are shown in Fig. 5. We observe that calibration performance is slightly higher on ImageNet-Adv when the hard calibration set is more difficult, while the performance on the other two datasets drops at the same time. Moreover, we find the optimal difficulty is different for various test sets. That said, by setting $d = 10$, we generally have good performance, and it is important to note that this

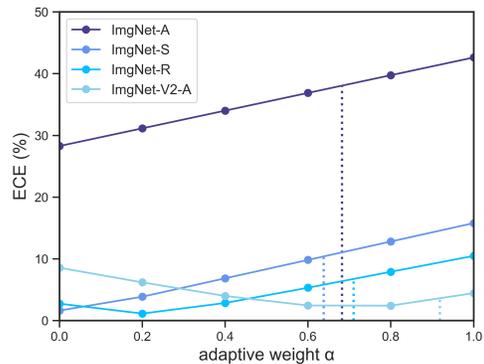


Figure 7. **Densely sampled values of α (0 to 1) vs. our computed α (Eq. 6)** Comparing with the densely sampled values of α , computed α (shown by the dotted vertical line) is close to the optimal value with reasonable difficulty for each test set.

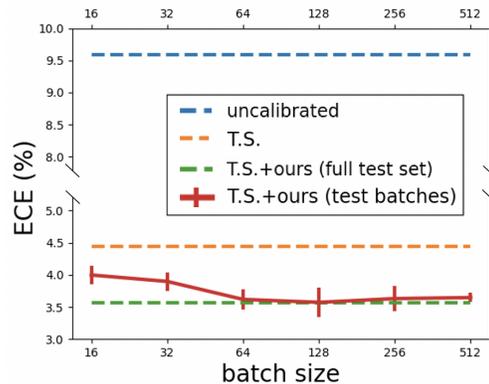


Figure 8. **Effectiveness of α computed in test batches of different sizes (16, 64, 128, 256 and 512).** Comparing with calculating α on the full set, using test batches yields similar ECE (%) especially when the batch size is at least 64. Temperature scaling (T.S.) is used as the baseline calibrator for our ACE. We also include the original baseline results in the figure (e.g., org T.S.).

difficulty level is considerably high (equivalent to 9.09% classification accuracy) and thus covers most test scenarios.

Comparing fixed weighting schemes with the adaptive weight We compare the adaptive weight ($\alpha = \frac{\text{avgConf}(\mathcal{D}_{test})}{\text{avgConf}(\mathcal{D}_o)}$) with setting α to a few fixed values 0.2, 0.5, and 0.9. The difficulty level for the out-of-distribution situation is 10. We evaluate the three calibration baselines on the six out-of-distribution test sets and one in-distribution test set using ResNet-152 as backbone and use the mean ECE (%) value over all the seven test sets (six OOD datasets and one ID test set) as evaluation metric.

As shown in Fig. 6, when applying our ACE on Spline, using $\alpha = 0.2$ and $\alpha = 0.5$ deteriorate calibration performance, while $\alpha = 0.9$ improves the baseline. However,

without test labels, it is infeasible to set an appropriate α for each test set. Moreover, the test sets are changed, setting fixed values of α might be effective in some cases and be less useful in others. In contrast, our designed test-set-adaptive α (Eq. 6) is shown to improve the baselines on various OOD test sets. Also, we report the value of α used for each test set in Table 1 and Table 2 in the supplemental material.

Optimal adaptive α vs. our computed α (Eq. 5). We compare both values in Fig. 7. First, for datasets with normal difficulty (e.g., ImageNet-V2-A), the value computed by our scheme is quite close to the optimal value. Second, for extremely difficult datasets such as ImageNet-S and ImageNet-A, α computed by our proposed method is less optimal. That said, we emphasize that in practice it is infeasible to do a greedy search because the images of test set are unlabeled, where our ACE method is generally useful.

Test data are given in batch. In real-world scenarios, test data may not all be accessible. Here we study how the calibration performance changes when test data are given in batches of various sizes. In Fig. 8, α is calculated from test batches of various sizes. We observe our method still achieves improvement over the temperature scaling baseline and has similar ECE with the method computed on the full test set under reasonably large batch sizes (≥ 64).

5. Conclusion

This paper studies how to calibrate a model on OOD datasets. Our important contribution is diagnosing why existing post-hoc algorithms fail on OOD test sets. Specifically, we report the difficulty of the calibration set influences the calibration function learning, and in other words, an OOD test set would witness poor calibration performance if the calibration set does not have an appropriate difficulty level. Realizing the importance of calibration set difficulty, we design a simple and effective method named adaptive calibrator ensemble (ACE) which combines the outputs of two calibrators trained on datasets with extreme difficulties. We also demonstrate how the ensemble scheme works for temperature scaling. We show that ACE improves three commonly used calibration methods on various OOD calibration benchmarks (e.g., ImageNet-C and CIFAR-10-C) without degrading ID calibration performance. In future work, we would like to further study how the domain gap and calibration set difficulty interact with each other and thereby improve OOD calibration.

Acknowledgement. We thank all anonymous reviewers and AC for their constructive comments and valuable suggestions. This work was supported in part by the National Natural Science Foundation of China (72174042), the ARC Discovery Early Career Researcher Award (DE200101283), and the ARC Discovery Project (DP210102801). Correspondence to: Yuli Zou. Yuli and Weijian contribute equally.

References

- [1] Mari-Liis Allikivi and Meelis Kull. Non-parametric bayesian isotonic calibration: Fighting over-confidence in binary classification. In *ECML/PKDD (2)*, pages 103–120, 2019.
- [2] 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.
- [3] Tianshuo Cong, Xinlei He, Yun Shen, and Yang Zhang. Test-time poisoning attacks against test-time adaptation models. *arXiv:2308.08505*, 2023.
- [4] 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.
- [5] Weijian Deng and Liang Zheng. Are labels necessary for classifier accuracy evaluation? In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2021.
- [6] 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. *arXiv preprint arXiv:2010.11929*, 2020.
- [7] Yunye Gong, Xiao Lin, Yi Yao, Thomas G Dietterich, Ajay Divakaran, and Melinda Gervasio. Confidence calibration for domain generalization under covariate shift. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8958–8967, 2021.
- [8] 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.
- [9] 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.
- [10] Kartik Gupta, Amir Rahimi, Thalayasingam Ajanthan, Thomas Mensink, Cristian Sminchisescu, and Richard Hartley. Calibration of neural networks using splines. In *International Conference on Learning Representations*, 2021.
- [11] Marton Havasi, Rodolphe Jenatton, Stanislav Fort, Jeremiah Zhe Liu, Jasper Snoek, Balaji Lakshminarayanan, Andrew Mingbo Dai, and Dustin Tran. Training independent subnetworks for robust prediction. In *ICLR*, 2021.
- [12] 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.
- [13] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical analysis of out-of-distribution generalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8340–8349, 2021.

- [14] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *Proceedings of the International Conference on Learning Representations*, 2019.
- [15] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial examples. *CVPR*, 2021.
- [16] Tom Joy, Francesco Pinto, Ser-Nam Lim, Philip HS Torr, and Puneet K Dokania. Sample-dependent adaptive temperature scaling for improved calibration. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 14919–14926, 2023.
- [17] Ranganath Krishnan and Omesh Tickoo. Improving model calibration with accuracy versus uncertainty optimization. *Advances in Neural Information Processing Systems*, 33:18237–18248, 2020.
- [18] Meelis Kull, Miquel Perello Nieto, Markus Kängsepp, Telmo Silva Filho, Hao Song, and Peter Flach. Beyond temperature scaling: Obtaining well-calibrated multi-class probabilities with dirichlet calibration. *Advances in neural information processing systems*, 32, 2019.
- [19] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. *Advances in neural information processing systems*, 30, 2017.
- [20] Guangrui Li, Guoliang Kang, Yi Zhu, Yunchao Wei, and Yi Yang. Domain consensus clustering for universal domain adaptation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021.
- [21] Jeremiah Liu, Zi Lin, Shreyas Padhy, Dustin Tran, Tania Bedrax Weiss, and Balaji Lakshminarayanan. Simple and principled uncertainty estimation with deterministic deep learning via distance awareness. *Advances in Neural Information Processing Systems*, 33:7498–7512, 2020.
- [22] Yishay Mansour, Mehryar Mohri, and Afshin Rostamizadeh. Domain adaptation: Learning bounds and algorithms. In *Proc. COLT*, 2009.
- [23] Mahdi Pakdaman Naeini and Gregory F Cooper. Binary classifier calibration using an ensemble of near isotonic regression models. In *2016 IEEE 16th International Conference on Data Mining (ICDM)*, pages 360–369. IEEE, 2016.
- [24] Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, David Sculley, Sebastian Nowozin, Joshua V Dillon, Balaji Lakshminarayanan, and Jasper Snoek. Can you trust your model’s uncertainty? evaluating predictive uncertainty under dataset shift. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, 2019.
- [25] Sumegha Premchandrar, Sandeep Madireddy, Sanket Jantre, and Prasanna Balaprakash. Unified probabilistic neural architecture and weight ensembling improves model robustness. *arXiv preprint arXiv:2210.04083*, 2022.
- [26] Amir Rahimi, Amirreza Shaban, Ching-An Cheng, Richard Hartley, and Byron Boots. Intra order-preserving functions for calibration of multi-class neural networks. *Advances in Neural Information Processing Systems*, 33:13456–13467, 2020.
- [27] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet classifiers generalize to im-
- agenet? In *International Conference on Machine Learning*, pages 5389–5400. PMLR, 2019.
- [28] Tiago Salvador, Vikram Voleti, Alexander Iannantuono, and Adam Oberman. Improved predictive uncertainty using corruption-based calibration. *stat*, 1050:7, 2021.
- [29] Rohan Taori, Achal Dave, Vaishaal Shankar, Nicholas Carlini, Benjamin Recht, and Ludwig Schmidt. Measuring robustness to natural distribution shifts in image classification. *Advances in Neural Information Processing Systems*, 33:18583–18599, 2020.
- [30] Junjiao Tian, Dylan Yung, Yen-Chang Hsu, and Zsolt Kira. A geometric perspective towards neural calibration via sensitivity decomposition. *Advances in Neural Information Processing Systems*, 34, 2021.
- [31] Christian Tomani, Sebastian Gruber, Muhammed Ebrar Erdem, Daniel Cremers, and Florian Buettner. Post-hoc uncertainty calibration for domain drift scenarios. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10124–10132, 2021.
- [32] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Herve Jegou. Training data-efficient image transformers and distillation through attention. In *International Conference on Machine Learning*, volume 139, pages 10347–10357, July 2021.
- [33] Weijie Tu, Weijian Deng, Tom Gedeon, and Zheng Liang. A bag-of-prototypes representation for dataset-level applications. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023.
- [34] Joost Van Amersfoort, Lewis Smith, Yee Whye Teh, and Yarin Gal. Uncertainty estimation using a single deep deterministic neural network. In *International conference on machine learning*, pages 9690–9700. PMLR, 2020.
- [35] Haohan Wang, Songwei Ge, Zachary Lipton, and Eric P Xing. Learning robust global representations by penalizing local predictive power. In *Advances in Neural Information Processing Systems*, pages 10506–10518, 2019.
- [36] Ximei Wang, Mingsheng Long, Jianmin Wang, and Michael Jordan. Transferable calibration with lower bias and variance in domain adaptation. *Advances in Neural Information Processing Systems*, 33:19212–19223, 2020.
- [37] Yeming Wen, Dustin Tran, and Jimmy Ba. Batchensemble: an alternative approach to efficient ensemble and lifelong learning. *arXiv preprint arXiv:2002.06715*, 2020.
- [38] Yeming Wen, Paul Vicol, Jimmy Ba, Dustin Tran, and Roger Grosse. Flipout: Efficient pseudo-independent weight perturbations on mini-batches. *arXiv preprint arXiv:1803.04386*, 2018.
- [39] Jonathan Wenger, Hedvig Kjellström, and Rudolph Triebel. Non-parametric calibration for classification. In *International Conference on Artificial Intelligence and Statistics*, pages 178–190. PMLR, 2020.
- [40] Bianca Zadrozny and Charles Elkan. Obtaining calibrated probability estimates from decision trees and naive bayesian classifiers. In *Icml*, volume 1, pages 609–616. Citeseer, 2001.
- [41] Jize Zhang, Bhavya Kailkhura, and T Yong-Jin Han. Mix-n-match: Ensemble and compositional methods for uncertainty calibration in deep learning. In *International conference on machine learning*, pages 11117–11128. PMLR, 2020.