IMPROVING THE POST-HOC CALIBRATION OF MOD-ERN NEURAL NETWORKS WITH PROBE SCALING

Anonymous authors

Paper under double-blind review

Abstract

We present "probe scaling": a post-hoc recipe for calibrating the predictions of modern neural networks. Our recipe is inspired by several lines of work, which demonstrate that early layers in the neural network learn general rules whereas later layers specialize. We show how such observations can be utilized in a post-hoc manner to calibrate the predictions of trained neural networks by injecting linear probes on the network's intermediate representations. Similar to temperature scaling, probe scaling neither retrains the architecture nor requires significantly more parameters. Unlike temperature scaling, however, it utilizes intermediate layers in the neural network. We demonstrate that probe scaling improves performance over temperature scaling on benchmark datasets across all five metrics: expected calibration error (ECE), negative log-likelihood, Brier score, classification accuracy, and the area under the ROC curve.

1 INTRODUCTION

Calibration is a measure of how accurate a model is in estimating its own uncertainty. It is an important metric for evaluating predictive models, particularly in critical domains such as medical diagnosis. In the latter case, having the ability to highlight predictions with low confidence is crucial in order for them to be handled separately, such as by referring to consultants for second opinions or by supplying additional information (e.g. lab tests), among others (Kompa et al., 2021). Providing an accurate estimate of uncertainty is also important in *selective predictions* (a.k.a. reject option classifiers) in which one optimizes a weighted combination of the cost of abstention and the cost of providing an incorrect prediction (Geifman & El-Yaniv, 2017)). It is also important when debiasing trained models via post-processing, where group-specific thresholding rules are adjusted to maximize accuracy subject to prescribed fairness guarantees (Corbett-Davies et al., 2017; Menon & Williamson, 2018; Celis et al., 2019; Alabdulmohsin & Lucic, 2021).

In order to achieve such objectives, however, trained machine learning models need to be *calibrated*, i.e. their predicted probabilities of outcomes should closely match with the observed empirical frequencies. More formally, writing \mathcal{X} for the instance space (e.g. the space of images), $\mathcal{Y} = \{0, 1, 2, \dots, K-1\}$ for the target set, and Δ_K for the probability simplex in \mathbb{R}^K , a probabilistic multiclass classifier $f : \mathcal{X} \to \Delta_K$ that outputs class probabilities for K classes is said to be calibrated if (Guo et al., 2017):

$$\forall y \in \mathcal{Y} : \ p(\mathbf{y} = y \,|\, f(\mathbf{x}) = q) = q_y, \tag{1}$$

where q_y is the y-th coordinate of the probability distribution $q \in \Delta_K$. In a slightly weaker notion, $f : \mathcal{X} \to \Delta_K$ is said to be *classwise* calibrated if (Kull et al., 2019):

$$\forall y \in \mathcal{Y}: \ p(\mathbf{y} = y \,| f_y(\mathbf{x}) = q_y) = q_y. \tag{2}$$

Unfortunately, several lines of work observed that modern neural networks lack such calibration (Guo et al., 2017; Thulasidasan et al., 2019; Ovadia et al., 2019; Kumar & Sarawagi, 2019; Kull et al., 2019), particularly under distribution shift (Ovadia et al., 2019). To remedy this problem, various approaches have been proposed including Monte Carlo dropout (Gal & Ghahramani, 2016), deep ensembles (Lakshminarayanan et al., 2016), training multiple independent subnetworks (Havasi et al., 2021), and augmentation (Thulasidasan et al., 2019). Often, such calibration methods inject *diversity* by generating multiple predictions for the same instance $\mathbf{x} \in \mathcal{X}$.



Figure 1: In probe scaling, d linear probes are injected into the intermediate representations of the neural network. Each probe is trained to maximize a proper score, such as the cross-entropy. Afterward, the logits of all probes are concatenated together to form a new representation. Finally, a classifier with d parameters (shown as β 's) is trained to output the calibrated scores using a 1D convolutional layer with strides d (equivalent to matrix multiplication) followed by softmax activations.

Besides diversity-based approaches, one particularly effective calibration rule is *temperature scaling* (Guo et al., 2017). It is an extension of the classical Platt scaling algorithm (Platt et al., 1999) to the multiclass setting. Temperature scaling rescales the neural networks' logits by a scalar $\tau \in \mathbb{R}^+$ chosen according to a separate validation dataset to optimize a *proper* scoring rule, where proper scores are those that guarantee to yield a perfectly calibrated classifier at its minimum at the infinite data limit (Gneiting & Raftery, 2007). Proper scores include, for example, the negative log-likelihood (a.k.a. cross-entropy loss) and the Brier score (a.k.a. squared loss). Despite the fact that temperature scaling is a post-hoc rule that neither requires changes to the training procedure nor adds significantly more parameters, it has been observed to yield competitive results in practice for in-distribution data (Guo et al., 2017; Ovadia et al., 2019; Kull et al., 2019).

In this work, we propose extending temperature scaling to "probe scaling", which is a post-hoc calibration method that utilizes the learned representations of *intermediate layers* in the neural network. An overview of probe scaling is presented in Figure 1. First, given a deep neural network that can be partitioned into d blocks (e.g. convolutional blocks in the standard ResNet and VGG architectures (He et al., 2016; Simonyan & Zisserman, 2015)), the *k*-th linear probe is trained to predict the final class $\mathbf{y} \in \mathcal{Y}$ given the representation learned at block *k* for the instance $\mathbf{x} \in \mathcal{X}$. In total, d - 1 linear probes are trained accordingly while the *d*-th probe is an identity mapping that preserves the logits of the original classifier. Finally, the logits of all probes are concatenated to form a new representation in the logit space with $|\mathcal{Y}| \times d$ features as shown in Figure 1, where $|\mathcal{Y}|$ is the number of classes. A final linear classifier is trained to optimize a proper score, such as the cross-entropy, on the new representation. Full pseudocode is provided in Algorithm 1.

Clearly, temperature scaling is a particular instance of probe scaling in which a single probe is used (i.e. d = 1). However, unlike temperature scaling, probe scaling often *improves* the accuracy of the original classifier besides improving its calibration. In fact, as discussed in Section 3, probe scaling performs better than temperature scaling on all considered five metrics: expected calibration error (ECE), negative log-likelihood, Brier score, classification accuracy, and the area under the ROC curve. The following proposition provides a formal argument for the advantage of probe scaling over temperature scaling when the aggregator is trained to optimize a proper, convex, Lipschitz scoring function, such as the logistic loss in the binary classification setting:

Proposition 1 Let $f : \mathbb{R}^d \times \mathcal{X} \times \mathcal{Y} \to \mathbb{R}$ be a scoring function that is convex and ρ -Lipschitz continuous on its first argument. For a given instance space \mathcal{X} and a target set \mathcal{Y} , denote by $\mathbf{S} \in$

 $(\mathcal{X} \times \mathcal{Y})^N \sim \mathcal{D}^N$ a hold-out validation dataset of size N. Let β_{probe} be the aggregator weights in probe scaling that minimize the regularized empirical loss on the hold-out sample:

$$\beta_{probe} = \arg\min_{\beta \in \mathbb{R}^d} \mathbb{E}_{\mathbf{x}, \mathbf{y} \in \mathbf{S}} \left[(\lambda/2) ||\beta||^2 + f(\beta, \mathbf{x}, \mathbf{y}) \right], \tag{3}$$

Then, the probability that the solution of temperature scaling β_{temp} performs better than probe scaling β_{probe} with respect to f by more than $\epsilon > 0$ over the random choice of the hold-out sample is bounded by:

$$p_{\mathbf{S}\sim\mathcal{D}^{N}}\left\{\mathbb{E}_{\mathbf{x},\mathbf{y}\in\mathcal{D}}\left[f(\beta_{temp},\mathbf{x},\mathbf{y})\right] \leq \mathbb{E}_{\mathbf{x},\mathbf{y}\in\mathcal{D}}\left[f(\beta_{probe},\mathbf{x},\mathbf{y})\right] - \epsilon\right\} \leq \frac{1}{\epsilon} \left(\lambda ||\beta^{\star}||^{2} + \frac{2\rho^{2}}{\lambda N}\right), \quad (4)$$

where β^* is the optimal probe scaling parameters at the limit of an infinite validation data size $N \to \infty$. Consequently, by setting $\lambda = 1/\sqrt{N}$, the probability bound $(||\beta^*||^2 + 2\rho^2)/(\epsilon\sqrt{N})$ can be made arbitrarily small for a sufficiently large N for any fixed $\epsilon > 0$.

In addition, one informal justification for why probe scaling improves calibration (and accuracy) is that early layers tend to learn general-purpose representations whereas later layers specialize. This has been observed in several works, such as when studying transfer learning (Raghu et al., 2019; Yosinski et al., 2014), memorization (Arpit et al., 2017; Cohen et al., 2018), and pretraining with random labels (Maennel et al., 2020). In particular, "easy" examples tend to be classified correctly by all layers of the neural network while difficult examples tend to be classified correctly by the last layers alone (Cohen et al., 2018; Baldock et al., 2021). By injecting linear probes on the intermediate representations of the neural network, the model improves its estimate of its own uncertainty by comparing the predictions at different layers.

2 CONTRIBUTION AND RELATED WORK

Measures of Calibration To evaluate calibration, several scores have been introduced in the literature. One popular score is the expected calibration error (ECE), which in its confidence-based version is given by (Naeini et al., 2015; Guo et al., 2017; Nixon et al., 2019):

$$ECE = \sum_{j=1}^{M} \frac{|B_j|}{n} |\operatorname{acc}(B_j) - p(B_j)|,$$
(5)

where $\{B_j\}_{j=1,...,M}$ are the bins, $|B_j|$ is the size of the bin, $acc(B_j)$ is the model's accuracy within the bin, and $p(B_j)$ is the model's average confidence. While often used as a default measure of calibration, ECE has its own limitations including high sensitivity to the number of chosen bins (Nixon et al., 2019). As a result, proper scoring rules are used as well. As discussed earlier, a scoring function is called proper if it is optimized by the true probability distribution, and it is called *strictly proper* if it has a unique optima (Gneiting & Raftery, 2007). One commonly used proper scoring function is the Brier score (Brier et al., 1950):

Brier =
$$\mathbb{E}_{\mathbf{x},\mathbf{y}} \Big[\sum_{y \in \mathcal{Y}} p(y \mid \mathbf{x})^2 - 2p(\mathbf{y} \mid \mathbf{x}) \Big],$$
 (6)

where **y** is the true label and $p(y | \mathbf{x})$ is the probability assigned by the model to the class $y \in \mathcal{Y}$ when the instance is **x**. Note that the Brier score is an affine transformation of the square loss and can be negative when Equation 6 is used. Besides, the negative log-likelihood (NLL), a.k.a. cross entropy, is a second proper score, which assigns a more aggressive penalty than the Brier score for events that occur when their predicted probabilities were quite small (Bickel, 2007). In our evaluation, we use the implementation of the ECE, Brier score, and NLL in Tensorflow (Abadi et al., 2015).

Calibration Methods. Several methods have been proposed to improve the calibration of neural networks. These include regularization-based methods, such as the Maximum Mean Calibration Error (MMCE), which adds a penalty term to the loss that is analogous to the Maximum Mean Discrepancy (MMD) method for testing the difference between two probability distributions (Kumar et al., 2018). They also include Bayesian methods, such as Stochastic Weight Average Gaussian (SWAG), which estimates uncertainty by sampling from the posterior distribution of the weights



Figure 2: Reliability diagrams are plotted for all five calibration methods when ResNet50 is trained on CIFAR100. Red shaded regions indicate the discrepancy between the confidence of the classifier and its accuracy. In general, both temperature scaling and probe scaling perform better than the other post-hoc methods with probe scaling performing best overall (see details in Section 3).

(Maddox et al., 2019). Other popular choices include ensemble-based approaches such as by training multiple networks with different seeds (Lakshminarayanan et al., 2017), augmentation-based methods such as Mix-Up (Thulasidasan et al., 2019), or variational inference based methods such as the Monte Carlo dropout (Gal & Ghahramani, 2016).

The proposed probe scaling method, on the other hand, is a *post-hoc* recipe for improving the calibration of modern neural networks. It extends temperature scaling (Guo et al., 2017) to intermediate layers in the neural network. Besides temperature scaling, an alternative post-hoc recipe for calibrating the predictions of neural networks is isotonic regression, which learns a monotonic transformation of the model's predictions into the true labels (Zadrozny & Elkan, 2002). In Gupta et al. (2021), another post-hoc recipe is proposed that is inspired by the Kolmogorov-Smirnov test, which approximates cumulative distribution functions using splines and estimates uncertainty using derivatives. For illustration, Figure 2 displays the reliability diagrams of these post-hoc rules when applied to ResNet50 (He et al., 2016) on the CIFAR100 dataset (Krizhevsky, 2009). A more thorough evaluation of all three post-hoc methods as well as probe scaling is provided in Section 3.

Statement of Contribution. In summary, we propose extending temperature scaling to intermediate layers in the neural network using linear probes. We provide a comprehensive evaluation showing that "probe scaling" outperforms other post-hoc rules in common benchmark datasets, including under distribution shift. We also demonstrate that probe scaling can complement other existing techniques, such as Monte Carlo dropout.

3 EXPERIMENTAL SETUP

3.1 COMPARISON WITH OTHER POST-HOC METHODS

Data Splits. To calibrate deep neural networks using post-hoc methods, we split the data into three parts: a training dataset (90% of standard training split), a hold-out dataset (10% of standard training split), and a test dataset (standard test split). In all methods, the original architecture is trained on the training dataset while the post-hoc rule is trained on the hold-out dataset. Final results are reported on the test data. In addition, the linear probes in the proposed probe scaling algorithm are trained on the training dataset as well. Nevertheless, we remark that the latter is by no means necessary since the probes can be applied to architectures that were pretrained on large datasets, such as ImageNet-21k (a superset of ILSVRC2012 that contains 21k classes (Deng et al., 2009)).

Datasets and Architectures. In our experiments, we train the models on CIFAR10/100 (Krizhevsky, 2009), Street View House Numbers (SVHN) (Netzer et al., 2011), and Fashion-MNIST (Xiao et al., 2017). We use five architectures: VGG16 and VGG19 (Simonyan & Zisserman, 2015), ResNet18 and ResNet50 (He et al., 2016), and Wide-ResNet-28-10 (Zagoruyko & Komodakis, 2016). We also conduct experiments on ImageNet ILSVRC2012 (Deng et al., 2009) using ResNet50 (He et al., 2016). To assess the impact of distribution shift, we evaluate ImageNet trained models on both clean images as well as the ImageNet-C corrupted images (Hendrycks & Dietterich, 2019).

Training Procedure. For reproducibility, we use a simple input training pipeline. We implement standard architectures with default ℓ_2 regularization of 10^{-4} . We train the model for 200 epochs on

Algorithm 1 Pseudocode of Probe Scaling

Input: (1) Trained neural network with identified sequence of $d - 1 \ge 0$ blocks; (2) Hold-out dataset.

Output: Calibrated model.

Training:

- 1: Insert d-1 linear probes, where probe k is trained to predict the target $\mathbf{y} \in \mathcal{Y}$ given the output (intermediate representation) of block k.
- 2: Concatenate the logits of all d 1 probes in addition to the logits of the original neural network, forming a new representation with $d \times |\mathcal{Y}|$ features. Let:

	$r_{1,1}$	$r_{2,1}$	•••	$r_{d,1}$
$R(\mathbf{v}) =$	$r_{1,2}$	$r_{2,2}$	•••	$r_{d,2}$
$n(\mathbf{x}) =$		•••	•••	
	$\lfloor r_{1, \mathcal{Y} }$	$r_{2, \mathcal{Y} }$	•••	$r_{d, \mathcal{Y} }$

be the new representation, where $r_{i,y}$ is the logit of the *i*-th probe for the class y.

3: Initialize $\beta_0 = (0, 0, ..., 1)^T \in \mathbb{R}^d$.

4: Train an aggregator using the projected stochastic gradient descent (SGD) that minimizes:

$$\mathbb{E}_{\mathbf{x},\mathbf{y}} \left| l \left(\text{softmax} \left(R(\mathbf{x}) \beta \right), \mathbf{y} \right) \right|$$

over the *d* parameters $\beta \in \mathbb{R}^d$ subject to the constraint $\beta \ge 0$, and starting from the initial value β_0 , for some proper scoring function *l* such as the cross-entropy. Here, $R(\mathbf{x}) \beta$ is matrix multiplication, which can be implemented using 1D convolutional layers with strides *d*.

CIFAR10/100, for 100 epochs on SVHN, and for 50 epochs on Fashion-MNIST, all with a batch size of 128. We use SGD with an initial learning rate of 0.1 and momentum of 0.9. We decay the learning rate with cosine learning rate decay (Loshchilov & Hutter, 2016). On the input pipeline, we standardize the input to have values in the unit interval [0, 1] and use only right-left random flipping as data augmentation. For ImageNet experiments, we use the pretrained ResNet50 architecture available in Keras (Chollet et al., 2015). We train on NVIDIA Tesla V100 GPUs for the medium-size datasets and train the probes in ImageNet on 32 TPUs. After training, we freeze the weights.

In probe scaling, each probe comprises of a fully connected classifier head layer with softmax activations trained using the cross-entropy loss. Before the classifier head, we flatten the intermediate representation that the probe is connected to, and we don't use any non-linearity between the backbone layers and the probe. For simplicity and reproducibility, we place a linear probe after every convolutional block in the ResNet and Wide-ResNet architectures, and after every convolutional layer in the VGG architectures. This results in 9 probes in ResNet18 (including the logits of the original classifier), 19 probes in ResNet50, 13 probes in Wide-ResNet-28-10, 17 probes in VGG16, and 20 probes in VGG19. Since the backbone network is frozen, training the probes can be carried out quickly. We train each probe with a learning rate of 0.005 and momentum 0.9 for 50 epochs, although we observe that they converge much faster. After each probe is trained, we also freeze it.

The final step in probe scaling is to train the final model, which aggregates the logits of all probes. It takes as input the concatenated output of each probe and performs a strided 1D convolution on it (See Figure 1 and Algorithm 1), which is equivalent to using one parameter per probe. We train the aggregator model for 50 epochs with a fixed learning rate of 0.005, using the hold-out part of the dataset. We repeat all experiments five times and report averages. Section 4.1 presents the results.

3.2 COMPLEMENTING EXISTING ALGORITHMS

The proposed probe scaling algorithm can also be combined with other existing techniques to improve calibration further. We illustrate this using the popular Monte Carlo dropout method (Gal & Ghahramani, 2016). In Monte Carlo dropout, multiple predictions are generated for the same instance $\mathbf{x} \in \mathcal{X}$ by adding dropout layers to the model and using them at *both* training and test time. Then, the model's uncertainty can be estimated by averaging the predictions.

In our experiments, we use the implementation of the Monte Carlo dropout available at the Uncertainty Baselines Benchmark (Nado et al., 2021). We use the same implementation details and

		$\text{ECE}\downarrow$	$\mathrm{NLL}\downarrow$	Brier \downarrow	ACC \uparrow	AUC \uparrow
vgg16	BASELINE TEMP ISOTONIC SPLINES PROBE	$\begin{array}{c} 0.201 \pm 0.006 \\ 0.041 \pm 0.003 \\ 0.069 \pm 0.005 \\ 0.049 \pm 0.010 \\ \textbf{0.025} \pm \textbf{0.004} \end{array}$	$\begin{array}{c} 1.759 \pm 0.044 \\ 1.238 \pm 0.026 \\ 1.392 \pm 0.039 \\ 3.382 \pm 0.231 \\ \textbf{0.981} \pm \textbf{0.020} \end{array}$	$\begin{array}{c} -0.513 \pm 0.011 \\ -0.582 \pm 0.008 \\ -0.565 \pm 0.009 \\ -0.273 \pm 0.030 \\ \textbf{-0.628 \pm 0.007} \end{array}$	$\begin{array}{c} 0.698 \pm 0.007 \\ 0.698 \pm 0.007 \\ 0.687 \pm 0.006 \\ 0.459 \pm 0.025 \\ \textbf{0.726} \pm \textbf{0.006} \end{array}$	$\begin{array}{c} 0.923 \pm 0.002 \\ 0.969 \pm 0.002 \\ 0.959 \pm 0.003 \\ 0.847 \pm 0.016 \\ \textbf{0.984} \pm \textbf{0.001} \end{array}$
vgg19	BASELINE TEMP ISOTONIC SPLINES PROBE	$\begin{array}{c} 0.217 \pm 0.005 \\ 0.046 \pm 0.005 \\ 0.057 \pm 0.007 \\ 0.052 \pm 0.012 \\ \textbf{0.026} \pm \textbf{0.003} \end{array}$	$\begin{array}{c} 2.002 \pm 0.040 \\ 1.303 \pm 0.022 \\ 1.436 \pm 0.029 \\ 3.929 \pm 0.313 \\ \textbf{0.969} \pm \textbf{0.018} \end{array}$	$\begin{array}{c} -0.497 \pm 0.010 \\ -0.581 \pm 0.008 \\ -0.568 \pm 0.008 \\ -0.221 \pm 0.033 \\ \textbf{-0.633} \pm 0.006 \end{array}$	$\begin{array}{c} 0.697 \pm 0.006 \\ 0.697 \pm 0.006 \\ 0.686 \pm 0.007 \\ 0.412 \pm 0.035 \\ \textbf{0.731} \pm \textbf{0.005} \end{array}$	$\begin{array}{c} 0.908 \pm 0.002 \\ 0.963 \pm 0.001 \\ 0.956 \pm 0.002 \\ 0.818 \pm 0.020 \\ \textbf{0.986 \pm 0.001} \end{array}$
RESNET18	BASELINE TEMP ISOTONIC SPLINES PROBE	$\begin{array}{c} 0.119 \pm 0.004 \\ 0.037 \pm 0.004 \\ 0.060 \pm 0.005 \\ 0.048 \pm 0.011 \\ \textbf{0.025} \pm \textbf{0.003} \end{array}$	$\begin{array}{c} 1.183 \pm 0.033 \\ 1.047 \pm 0.025 \\ 1.245 \pm 0.031 \\ 2.975 \pm 0.285 \\ \textbf{0.925} \pm \textbf{0.021} \end{array}$	$\begin{array}{c} -0.604 \pm 0.009 \\ -0.627 \pm 0.007 \\ -0.597 \pm 0.008 \\ -0.328 \pm 0.035 \\ \textbf{-0.651} \pm \textbf{0.007} \end{array}$	$\begin{array}{c} 0.728 \pm 0.006 \\ 0.728 \pm 0.006 \\ 0.706 \pm 0.006 \\ 0.486 \pm 0.029 \\ \textbf{0.746} \pm \textbf{0.005} \end{array}$	$\begin{array}{c} 0.960 \pm 0.002 \\ 0.979 \pm 0.001 \\ 0.968 \pm 0.002 \\ 0.870 \pm 0.020 \\ \textbf{0.985 \pm 0.001} \end{array}$
resnet50	BASELINE TEMP ISOTONIC SPLINES PROBE	$\begin{array}{c} 0.185 \pm 0.006 \\ 0.021 \pm 0.003 \\ 0.058 \pm 0.006 \\ 0.041 \pm 0.007 \\ \textbf{0.018} \pm \textbf{0.003} \end{array}$	$\begin{array}{c} 1.675 \pm 0.033 \\ 0.965 \pm 0.015 \\ 1.353 \pm 0.028 \\ 2.769 \pm 0.202 \\ \textbf{0.949} \pm \textbf{0.015} \end{array}$	$\begin{array}{c} -0.561 \pm 0.009 \\ -0.634 \pm 0.006 \\ -0.545 \pm 0.008 \\ -0.343 \pm 0.025 \\ \textbf{-0.637} \pm \textbf{0.006} \end{array}$	$\begin{array}{c} 0.731 \pm 0.006 \\ 0.731 \pm 0.006 \\ 0.670 \pm 0.007 \\ 0.499 \pm 0.021 \\ \textbf{0.733} \pm \textbf{0.007} \end{array}$	$\begin{array}{c} 0.933 \pm 0.002 \\ 0.983 \pm 0.001 \\ 0.967 \pm 0.002 \\ 0.885 \pm 0.014 \\ \textbf{0.984} \pm \textbf{0.001} \end{array}$
WIDE-RESNET-28-10	BASELINE TEMP ISOTONIC SPLINES PROBE	$\begin{array}{c} 0.054 \pm 0.004 \\ 0.054 \pm 0.004 \\ \textbf{0.046} \pm \textbf{0.004} \\ 0.054 \pm 0.011 \\ 0.053 \pm 0.004 \end{array}$	$\begin{array}{c} 0.890 \pm 0.017 \\ 0.890 \pm 0.017 \\ 0.980 \pm 0.024 \\ 3.312 \pm 0.263 \\ \textbf{0.887} \pm \textbf{0.017} \end{array}$	-0.691 ± 0.005 -0.691 ± 0.005 -0.686 ± 0.006 -0.300 ± 0.035 -0.691 ± 0.005	$\begin{array}{c} \textbf{0.785} \pm \textbf{0.004} \\ \textbf{0.785} \pm \textbf{0.004} \\ \textbf{0.775} \pm \textbf{0.005} \\ \textbf{0.470} \pm \textbf{0.029} \\ \textbf{0.785} \pm \textbf{0.004} \end{array}$	$\begin{array}{c} 0.981 \pm 0.001 \\ 0.981 \pm 0.001 \\ 0.972 \pm 0.002 \\ 0.839 \pm 0.020 \\ \textbf{0.982 \pm 0.001} \end{array}$

Table 1: Empirical evaluation of post-hoc calibration methods on CIFAR100 (Krizhevsky, 2009) for the five metrics: ECE, negative log-likelihood (NLL), Brier score, accuracy (ACC), and AUC.

Table 2: This table lists the *p*-values using the non-parametric Wilcoxon signed-rank test for each post-hoc method compared against probe scaling. In all cases, statistics are in favor of probe scaling. The *p*-values shown in **bold** remain statistically significant at the 95% confidence level even after correcting for multiple hypothesis testing using Holm's step-down procedure (Demšar, 2006).

	ECE	NLL	Brier	ACC	AUC
BASELINE TEMP ISOTONIC SPLINES	$\begin{array}{c} {\bf 1} \times {\bf 10^{-4}} \\ {\bf 8} \times {\bf 10^{-3}} \\ {\bf 3} \times {\bf 10^{-1}} \\ {\bf 1} \times {\bf 10^{-1}} \end{array}$	$egin{array}{c} 8 imes 10^{-5} \ 3 imes 10^{-4} \ 3 imes 10^{-4} \ 9 imes 10^{-5} \end{array}$	$\begin{array}{c} 9\times 10^{-5} \\ 9\times 10^{-5} \\ 5\times 10^{-4} \\ 1\times 10^{-4} \end{array}$	$egin{array}{c} 2 imes 10^{-4} \ 2 imes 10^{-4} \ 1 imes 10^{-3} \ 1 imes 10^{-4} \end{array}$	$egin{array}{c} 9 imes 10^{-5} \ 9 imes 10^{-4} \ 2 imes 10^{-4} \ 9 imes 10^{-5} \end{array}$

hyper-parameters suggested by the benchmark, which uses the Wide-ResNet-28-10 architecture, and follow it by probe scaling as described earlier in Section 3.1. We also test on different values of the dropout rate. The final results are presented in Section 4.2.

4 RESULTS AND DISCUSSION

4.1 POST-HOC CALIBERATION RESULTS

Tables 1-5 present the empirical results of using probe scaling to calibrate neural networks in a posthoc manner. Five metrics are used in this evaluation: expected calibration error (ECE), negative loglikelihood (NLL), Brier score (Brier), accuracy (ACC), and the area under the ROC curve (AUC). We compare probe scaling to vanilla training (baseline), temperature scaling (Guo et al., 2017), isotonic regression (Zadrozny & Elkan, 2002), and splines (Gupta et al., 2021). We observe that probe scaling consistently outperforms other post-hoc methods in all metrics except on ECE in which no method seems to dominate. Also, probe scaling is the only post-hoc method that consistently improves the accuracy of the model and its AUC. For ImageNet and ImageNet-C, on the other hand, probe scaling seems to offer a modest gain over temperature scaling as shown in Table 6.

Statistical Significance. To determine if the improvement of probe scaling over other post-hoc methods is statistically significant, we used the non-parametric Wilcoxon signed-rank test and cor-

		ECE \downarrow	NLL \downarrow	Brier \downarrow	ACC \uparrow	AUC \uparrow
vgg16	BASELINE TEMP ISOTONIC SPLINES PROBE	$\begin{array}{c} 0.060 \pm 0.002 \\ 0.025 \pm 0.003 \\ 0.018 \pm 0.003 \\ \textbf{0.013} \pm \textbf{0.002} \\ 0.016 \pm 0.003 \end{array}$	$\begin{array}{c} 0.423 \pm 0.018 \\ 0.263 \pm 0.009 \\ 0.260 \pm 0.009 \\ 0.298 \pm 0.011 \\ \textbf{0.226} \pm \textbf{0.007} \end{array}$	$\begin{array}{c} -0.869 \pm 0.004 \\ -0.885 \pm 0.004 \\ -0.887 \pm 0.003 \\ -0.886 \pm 0.003 \\ \textbf{-0.892} \pm 0.004 \end{array}$	$\begin{array}{c} 0.926 \pm 0.002 \\ 0.926 \pm 0.002 \\ 0.924 \pm 0.003 \\ 0.924 \pm 0.003 \\ \textbf{0.927} \pm \textbf{0.002} \end{array}$	$\begin{array}{c} 0.979 \pm 0.001 \\ 0.994 \pm 0.000 \\ 0.993 \pm 0.001 \\ 0.992 \pm 0.001 \\ \textbf{0.996} \pm \textbf{0.000} \end{array}$
vgg19	BASELINE TEMP ISOTONIC SPLINES PROBE	$\begin{array}{c} 0.061 \pm 0.003 \\ 0.023 \pm 0.003 \\ 0.017 \pm 0.004 \\ \textbf{0.012} \pm \textbf{0.002} \\ 0.017 \pm 0.002 \end{array}$	$\begin{array}{c} 0.444 \pm 0.021 \\ 0.269 \pm 0.009 \\ 0.269 \pm 0.013 \\ 0.313 \pm 0.018 \\ \textbf{0.226} \pm \textbf{0.006} \end{array}$	$\begin{array}{c} -0.869 \pm 0.005 \\ -0.885 \pm 0.004 \\ -0.886 \pm 0.004 \\ -0.885 \pm 0.003 \\ \textbf{-0.892} \pm 0.003 \end{array}$	$\begin{array}{c} 0.927 \pm 0.002 \\ 0.927 \pm 0.002 \\ 0.926 \pm 0.003 \\ 0.924 \pm 0.003 \\ \textbf{0.928} \pm \textbf{0.002} \end{array}$	$\begin{array}{c} 0.977 \pm 0.001 \\ 0.994 \pm 0.001 \\ 0.993 \pm 0.001 \\ 0.991 \pm 0.001 \\ \textbf{0.996} \pm \textbf{0.000} \end{array}$
RESNET18	BASELINE TEMP ISOTONIC SPLINES PROBE	$\begin{array}{c} 0.038 \pm 0.002 \\ 0.013 \pm 0.003 \\ 0.015 \pm 0.004 \\ \textbf{0.011} \pm \textbf{0.004} \\ 0.015 \pm 0.002 \end{array}$	$\begin{array}{c} 0.265 \pm 0.016 \\ 0.191 \pm 0.010 \\ 0.240 \pm 0.026 \\ 0.247 \pm 0.036 \\ \textbf{0.184 \pm 0.009} \end{array}$	$\begin{array}{c} -0.904 \pm 0.004 \\ \textbf{-0.912} \pm \textbf{0.004} \\ -0.893 \pm 0.011 \\ -0.903 \pm 0.014 \\ \textbf{-0.912} \pm \textbf{0.004} \end{array}$	$\begin{array}{c} \textbf{0.941} \pm \textbf{0.003} \\ \textbf{0.941} \pm \textbf{0.003} \\ \textbf{0.928} \pm \textbf{0.008} \\ \textbf{0.933} \pm \textbf{0.014} \\ \textbf{0.941} \pm \textbf{0.002} \end{array}$	$\begin{array}{c} 0.989 \pm 0.001 \\ \textbf{0.996} \pm \textbf{0.001} \\ 0.994 \pm 0.001 \\ 0.994 \pm 0.001 \\ \textbf{0.996} \pm \textbf{0.000} \end{array}$
resnet50	BASELINE TEMP ISOTONIC SPLINES PROBE	$\begin{array}{c} 0.046 \pm 0.002 \\ \textbf{0.010} \pm \textbf{0.002} \\ 0.028 \pm 0.008 \\ 0.046 \pm 0.016 \\ \textbf{0.010} \pm \textbf{0.002} \end{array}$	$\begin{array}{c} 0.341 \pm 0.016 \\ 0.190 \pm 0.006 \\ 0.440 \pm 0.065 \\ 0.740 \pm 0.146 \\ \textbf{0.189} \pm \textbf{0.006} \end{array}$	-0.893 ± 0.003 -0.907 ± 0.002 -0.791 ± 0.031 -0.718 ± 0.049 -0.907 ± 0.002	$\begin{array}{c} \textbf{0.936} \pm \textbf{0.002} \\ \textbf{0.936} \pm \textbf{0.002} \\ \textbf{0.857} \pm \textbf{0.022} \\ \textbf{0.769} \pm \textbf{0.041} \\ \textbf{0.936} \pm \textbf{0.002} \end{array}$	$\begin{array}{c} 0.986 \pm 0.001 \\ \textbf{0.997} \pm \textbf{0.000} \\ 0.986 \pm 0.003 \\ 0.961 \pm 0.014 \\ \textbf{0.997} \pm \textbf{0.000} \end{array}$
wide-resnet-28-10	BASELINE TEMP ISOTONIC SPLINES PROBE	$\begin{array}{c} 0.026 \pm 0.002 \\ 0.012 \pm 0.002 \\ 0.012 \pm 0.002 \\ \textbf{0.009} \pm \textbf{0.002} \\ 0.011 \pm 0.001 \end{array}$	$\begin{array}{c} 0.176 \pm 0.010 \\ 0.149 \pm 0.007 \\ 0.157 \pm 0.008 \\ 0.192 \pm 0.014 \\ \textbf{0.147} \pm \textbf{0.006} \end{array}$	$\begin{array}{c} -0.928 \pm 0.003 \\ \textbf{-0.932} \pm \textbf{0.003} \\ \textbf{-0.932} \pm \textbf{0.003} \\ -0.929 \pm 0.003 \\ \textbf{-0.932} \pm \textbf{0.003} \end{array}$	$\begin{array}{c} \textbf{0.956} \pm \textbf{0.002} \\ \textbf{0.956} \pm \textbf{0.002} \\ \textbf{0.955} \pm \textbf{0.002} \\ \textbf{0.953} \pm \textbf{0.003} \\ \textbf{0.956} \pm \textbf{0.002} \end{array}$	$\begin{array}{c} 0.994 \pm 0.001 \\ \textbf{0.997} \pm \textbf{0.000} \\ \textbf{0.997} \pm \textbf{0.001} \\ 0.995 \pm 0.001 \\ \textbf{0.997} \pm \textbf{0.000} \end{array}$

Table 3: Empirical evaluation of post-hoc calibration methods on CIFAR10 (Krizhevsky, 2009) for the five metrics: ECE, negative log-likelihood (NLL), Brier score, accuracy (ACC) and AUC.

Table 4: Empirical evaluation of post-hoc calibration methods on SVHN (Netzer et al., 2011) for the five metrics: ECE, negative log-likelihood (NLL), Brier score, accuracy (ACC) and AUC.

		ECE \downarrow	$\mathrm{NLL}\downarrow$	Brier \downarrow	ACC \uparrow	AUC ↑
vgg16	BASELINE TEMP ISOTONIC SPLINES PROBE	$\begin{array}{c} 0.021 \pm 0.003 \\ 0.012 \pm 0.002 \\ \textbf{0.006} \pm \textbf{0.001} \\ 0.013 \pm 0.004 \\ 0.011 \pm 0.002 \end{array}$	$\begin{array}{c} 0.227 \pm 0.009 \\ 0.218 \pm 0.011 \\ 0.208 \pm 0.009 \\ 0.420 \pm 0.033 \\ \textbf{0.191} \pm \textbf{0.009} \end{array}$	$\begin{array}{c} -0.907 \pm 0.005 \\ -0.909 \pm 0.005 \\ -0.914 \pm 0.004 \\ -0.869 \pm 0.009 \\ \textbf{-0.920} \pm 0.004 \end{array}$	$\begin{array}{c} 0.940 \pm 0.004 \\ 0.940 \pm 0.004 \\ 0.944 \pm 0.003 \\ 0.917 \pm 0.005 \\ \textbf{0.947} \pm \textbf{0.003} \end{array}$	$\begin{array}{c} 0.993 \pm 0.001 \\ 0.995 \pm 0.000 \\ 0.994 \pm 0.001 \\ 0.980 \pm 0.002 \\ \textbf{0.996} \pm \textbf{0.000} \end{array}$
vgg19	BASELINE TEMP ISOTONIC SPLINES PROBE	$\begin{array}{c} 0.019 \pm 0.002 \\ 0.012 \pm 0.002 \\ \textbf{0.006} \pm \textbf{0.001} \\ 0.013 \pm 0.004 \\ 0.010 \pm 0.003 \end{array}$	$\begin{array}{c} 0.205 \pm 0.011 \\ 0.199 \pm 0.009 \\ 0.194 \pm 0.010 \\ 0.408 \pm 0.044 \\ \textbf{0.181} \pm \textbf{0.007} \end{array}$	$\begin{array}{c} -0.918 \pm 0.004 \\ -0.920 \pm 0.004 \\ -0.922 \pm 0.004 \\ -0.877 \pm 0.010 \\ \textbf{-0.925} \pm 0.003 \end{array}$	$\begin{array}{c} 0.947 \pm 0.003 \\ 0.947 \pm 0.003 \\ 0.949 \pm 0.003 \\ 0.923 \pm 0.005 \\ \textbf{0.951} \pm \textbf{0.002} \end{array}$	$\begin{array}{c} 0.993 \pm 0.001 \\ 0.995 \pm 0.000 \\ 0.994 \pm 0.000 \\ 0.981 \pm 0.003 \\ \textbf{0.996} \pm \textbf{0.000} \end{array}$
RESNET18	BASELINE TEMP ISOTONIC SPLINES PROBE	$\begin{array}{c} 0.026 \pm 0.002 \\ 0.012 \pm 0.003 \\ \textbf{0.007} \pm \textbf{0.001} \\ 0.011 \pm 0.003 \\ 0.013 \pm 0.002 \end{array}$	$\begin{array}{c} 0.250 \pm 0.018 \\ 0.228 \pm 0.016 \\ 0.218 \pm 0.014 \\ 0.409 \pm 0.031 \\ \textbf{0.205 \pm 0.013} \end{array}$	$\begin{array}{c} -0.901 \pm 0.007 \\ -0.904 \pm 0.007 \\ -0.908 \pm 0.006 \\ -0.868 \pm 0.006 \\ \textbf{-0.912} \pm 0.006 \end{array}$	$\begin{array}{c} 0.937 \pm 0.005 \\ 0.937 \pm 0.005 \\ 0.940 \pm 0.004 \\ 0.917 \pm 0.004 \\ \textbf{0.942} \pm \textbf{0.004} \end{array}$	$\begin{array}{c} 0.991 \pm 0.001 \\ 0.995 \pm 0.001 \\ 0.994 \pm 0.001 \\ 0.981 \pm 0.003 \\ \textbf{0.996} \pm \textbf{0.001} \end{array}$
resnet50	BASELINE TEMP ISOTONIC SPLINES PROBE	$\begin{array}{c} 0.013 \pm 0.002 \\ 0.016 \pm 0.002 \\ \textbf{0.006} \pm \textbf{0.001} \\ 0.011 \pm 0.003 \\ 0.014 \pm 0.002 \end{array}$	$\begin{array}{c} 0.209 \pm 0.018 \\ 0.205 \pm 0.018 \\ \textbf{0.183} \pm \textbf{0.013} \\ 0.397 \pm 0.029 \\ 0.191 \pm 0.014 \end{array}$	$\begin{array}{c} -0.914 \pm 0.008 \\ -0.914 \pm 0.008 \\ \textbf{-0.924} \pm \textbf{0.006} \\ -0.879 \pm 0.008 \\ -0.919 \pm 0.007 \end{array}$	$\begin{array}{c} 0.944 \pm 0.006 \\ 0.944 \pm 0.006 \\ \textbf{0.952} \pm \textbf{0.004} \\ 0.925 \pm 0.005 \\ 0.947 \pm 0.005 \end{array}$	$\begin{array}{c} 0.994 \pm 0.001 \\ 0.995 \pm 0.001 \\ 0.995 \pm 0.001 \\ 0.981 \pm 0.002 \\ \textbf{0.996} \pm \textbf{0.000} \end{array}$
WIDE-RESNET-28-10	BASELINE TEMP ISOTONIC SPLINES PROBE	$\begin{array}{c} 0.020 \pm 0.002 \\ 0.009 \pm 0.001 \\ \textbf{0.005} \pm \textbf{0.001} \\ 0.012 \pm 0.004 \\ 0.009 \pm 0.001 \end{array}$	$\begin{array}{c} 0.196 \pm 0.010 \\ 0.181 \pm 0.011 \\ 0.175 \pm 0.009 \\ 0.421 \pm 0.042 \\ \textbf{0.174 \pm 0.008} \end{array}$	$\begin{array}{c} -0.923 \pm 0.005 \\ -0.925 \pm 0.005 \\ \textbf{-0.928} \pm \textbf{0.004} \\ -0.878 \pm 0.008 \\ -0.927 \pm 0.004 \end{array}$	$\begin{array}{c} 0.951 \pm 0.004 \\ 0.951 \pm 0.004 \\ \textbf{0.954} \pm \textbf{0.003} \\ 0.924 \pm 0.005 \\ 0.953 \pm 0.003 \end{array}$	$\begin{array}{c} 0.993 \pm 0.000 \\ \textbf{0.996} \pm \textbf{0.000} \\ 0.995 \pm 0.000 \\ 0.980 \pm 0.003 \\ \textbf{0.996} \pm \textbf{0.000} \end{array}$

rected it for multiple hypothesis testing using Holm's step down procedure (Demšar, 2006). Our analysis reveals that probe scaling outperforms both vanilla training (baseline) and temperature scaling in ECE, and outperforms all post-hoc methods on all remaining four metrics: NLL, Brier, accuracy, and AUC with a statistically significant evidence at the 95% confidence level. Details are provided in Table 2.

		$\text{ECE}\downarrow$	$\mathrm{NLL}\downarrow$	Brier \downarrow	ACC \uparrow	AUC ↑
vgg16	BASELINE TEMP ISOTONIC SPLINES PROBE	$\begin{array}{c} 0.049 \pm 0.007 \\ 0.017 \pm 0.003 \\ 0.016 \pm 0.003 \\ 0.016 \pm 0.005 \\ \textbf{0.010} \pm \textbf{0.002} \end{array}$	$\begin{array}{c} 0.323 \pm 0.039 \\ 0.274 \pm 0.021 \\ 0.273 \pm 0.018 \\ 0.340 \pm 0.077 \\ \textbf{0.199} \pm \textbf{0.007} \end{array}$	$\begin{array}{c} -0.858 \pm 0.010 \\ -0.866 \pm 0.009 \\ -0.872 \pm 0.008 \\ -0.857 \pm 0.028 \\ \textbf{-0.898} \pm \textbf{0.004} \end{array}$	$\begin{array}{c} 0.910 \pm 0.005 \\ 0.910 \pm 0.005 \\ 0.914 \pm 0.005 \\ 0.895 \pm 0.031 \\ \textbf{0.931} \pm \textbf{0.003} \end{array}$	$\begin{array}{c} 0.988 \pm 0.003 \\ 0.994 \pm 0.001 \\ 0.993 \pm 0.001 \\ 0.991 \pm 0.003 \\ \textbf{0.997} \pm \textbf{0.000} \end{array}$
vgg19	BASELINE TEMP ISOTONIC SPLINES PROBE	$\begin{array}{c} 0.040 \pm 0.012 \\ 0.015 \pm 0.003 \\ 0.016 \pm 0.003 \\ \textbf{0.014} \pm \textbf{0.004} \\ 0.023 \pm 0.025 \end{array}$	$\begin{array}{c} 0.298 \pm 0.064 \\ \textbf{0.267} \pm \textbf{0.048} \\ 0.269 \pm 0.028 \\ 0.320 \pm 0.070 \\ 0.280 \pm 0.178 \end{array}$	$\begin{array}{c} -0.864 \pm 0.026 \\ -0.869 \pm 0.023 \\ -0.873 \pm 0.014 \\ -0.862 \pm 0.030 \\ \textbf{-0.889} \pm \textbf{0.026} \end{array}$	$\begin{array}{c} 0.912 \pm 0.015 \\ 0.912 \pm 0.015 \\ 0.915 \pm 0.010 \\ 0.902 \pm 0.031 \\ \textbf{0.928} \pm \textbf{0.010} \end{array}$	$\begin{array}{c} 0.990 \pm 0.003 \\ \textbf{0.994} \pm \textbf{0.002} \\ \textbf{0.994} \pm \textbf{0.001} \\ 0.992 \pm 0.003 \\ 0.992 \pm 0.010 \end{array}$
RESNET18	BASELINE TEMP ISOTONIC SPLINES PROBE	$\begin{array}{c} 0.053 \pm 0.006 \\ \textbf{0.014} \pm \textbf{0.003} \\ 0.020 \pm 0.005 \\ 0.030 \pm 0.015 \\ \textbf{0.014} \pm \textbf{0.004} \end{array}$	$\begin{array}{c} 0.358 \pm 0.043 \\ 0.261 \pm 0.019 \\ 0.323 \pm 0.030 \\ 0.568 \pm 0.200 \\ \textbf{0.210} \pm \textbf{0.005} \end{array}$	-0.863 ± 0.011 -0.875 ± 0.009 -0.847 ± 0.016 -0.773 ± 0.072 -0.897 ± 0.003	$\begin{array}{c} 0.915 \pm 0.006 \\ 0.915 \pm 0.006 \\ 0.898 \pm 0.010 \\ 0.812 \pm 0.063 \\ \textbf{0.931} \pm \textbf{0.002} \end{array}$	$\begin{array}{c} 0.986 \pm 0.002 \\ 0.994 \pm 0.001 \\ 0.991 \pm 0.001 \\ 0.975 \pm 0.017 \\ \textbf{0.996} \pm \textbf{0.000} \end{array}$
resnet50	BASELINE TEMP ISOTONIC SPLINES PROBE	$\begin{array}{c} 0.064 \pm 0.006 \\ 0.014 \pm 0.002 \\ 0.017 \pm 0.004 \\ 0.026 \pm 0.014 \\ \textbf{0.012} \pm \textbf{0.002} \end{array}$	$\begin{array}{c} 0.433 \pm 0.021 \\ 0.331 \pm 0.012 \\ 0.318 \pm 0.017 \\ 0.480 \pm 0.160 \\ \textbf{0.282 \pm 0.021} \end{array}$	$\begin{array}{c} -0.823 \pm 0.008 \\ -0.836 \pm 0.007 \\ -0.845 \pm 0.008 \\ -0.801 \pm 0.057 \\ \textbf{-0.859} \pm 0.010 \end{array}$	$\begin{array}{c} 0.887 \pm 0.004 \\ 0.887 \pm 0.004 \\ 0.896 \pm 0.004 \\ 0.849 \pm 0.054 \\ \textbf{0.904} \pm \textbf{0.006} \end{array}$	$\begin{array}{c} 0.984 \pm 0.001 \\ 0.992 \pm 0.001 \\ 0.992 \pm 0.001 \\ 0.982 \pm 0.012 \\ \textbf{0.994} \pm \textbf{0.001} \end{array}$
wide-resnet-28-10	BASELINE TEMP ISOTONIC SPLINES PROBE	$\begin{array}{c} 0.043 \pm 0.004 \\ 0.013 \pm 0.002 \\ 0.013 \pm 0.002 \\ 0.013 \pm 0.002 \\ \textbf{0.010} \pm \textbf{0.001} \end{array}$	$\begin{array}{c} 0.300 \pm 0.014 \\ 0.252 \pm 0.009 \\ \textbf{0.242} \pm \textbf{0.008} \\ 0.270 \pm 0.011 \\ \textbf{0.242} \pm \textbf{0.007} \end{array}$	$\begin{array}{c} -0.873 \pm 0.005 \\ -0.879 \pm 0.004 \\ \textbf{-0.886} \pm 0.003 \\ -0.884 \pm 0.003 \\ -0.882 \pm 0.003 \end{array}$	$\begin{array}{c} 0.920 \pm 0.003 \\ 0.920 \pm 0.003 \\ \textbf{0.925} \pm \textbf{0.003} \\ 0.923 \pm 0.003 \\ 0.922 \pm 0.002 \end{array}$	$\begin{array}{c} 0.989 \pm 0.001 \\ 0.994 \pm 0.001 \\ 0.994 \pm 0.001 \\ 0.994 \pm 0.001 \\ \textbf{0.995 \pm 0.000} \end{array}$

Table 5: Empirical evaluation of post-hoc calibration methods on Fashion MNIST (Xiao et al., 2017) for the five metrics: ECE, negative log-likelihood (NLL), Brier score, accuracy (ACC) and AUC.

Table 6: Empirical evaluation on ResNet50 trained on ImageNet ILSVRC2012 (Deng et al., 2009) and ImageNet-C (Hendrycks & Dietterich, 2019). Here, the splines-based method is excluded from the comparison because it failed to perform better than the baseline.

		$ECE\downarrow$	$\mathrm{NLL}\downarrow$	Brier \downarrow	ACC \uparrow	AUC ↑
ImageNet	BASELINE	0.0640	1.1481	-0.6173	0.7276	0.6659
	ISOTONIC	0.0447	1.3849	-0.5978	0.7093	0.6947
	TEMP	0.0247	1.1066	-0.6220	0.7276	0.6735
	PROBE	0.0239	1.1065	-0.6220	0.7272	0.6736
ImageNet-C	BASELINE	0.0846	1.6544	-0.5005	0.6316	0.6548
	ISOTONIC	0.0402	1.8061	-0.4965	0.6249	0.6858
	TEMP	0.0271	1.5903	-0.5102	0.6316	0.6646
	PROBE	0.0259	1.5903	-0.5101	0.6316	0.6646

4.2 COMPLEMENTING METHODS RESULTS

Table 7 presents the empirical results when combining Monte Carlo dropout with probe scaling. We evaluate on the Uncertainty Baselines Benchmark (Nado et al., 2021) as mentioned earlier, which uses the Wide-ResNet-28-10 architecture. The results indicate that in most of the cases, probe scaling improves the calibration and uncertainty estimation of the model when combined with Monte Carlo dropout, sometimes quite significantly such as in the Fashion MNIST dataset. Out of the four datasets, only SVHN does not reveal a notable gain by combining probe scaling with Monte Carlo dropout.

4.3 CONTRIBUTION OF PROBES.

In Figure 3, we plot the contribution of the probes as a function of their relative depths for models trained on CIFAR100. One stark difference appears between the ResNet architectures with skip connections and the classical VGG architectures. In VGG16/19, probes assigned to the *early* layers in the neural network play a significant role in estimating the model's uncertainty. In ResNet architectures, on the other hand, only the probes assigned to the upper layers play a significant role.

			ECE \downarrow	$\mathrm{NLL}\downarrow$	Brier \downarrow	ACC \uparrow	AUC \uparrow		
CIFAR10									
Dropout rate	0.1	MC Dropout Probes + MC Dropout	$\begin{array}{c} 0.023 \pm 0.001 \\ \textbf{0.010} \pm \textbf{0.001} \end{array}$	$\begin{array}{c} 0.160 \pm 0.004 \\ \textbf{0.144} \pm \textbf{0.003} \end{array}$	-0.932 ± 0.001 -0.935 ± 0.001	$\begin{array}{c} 0.958 \pm 0.001 \\ 0.957 \pm 0.001 \end{array}$	$\begin{array}{c} 0.995 \pm 0.000 \\ \textbf{0.997} \pm \textbf{0.000} \end{array}$		
Dropout rate	0.2	MC Dropout Probes + MC Dropout	$\begin{array}{c} 0.019 \pm 0.001 \\ \textbf{0.005} \pm \textbf{0.000} \end{array}$	$\begin{array}{c} 0.151 \pm 0.002 \\ \textbf{0.144} \pm \textbf{0.001} \end{array}$	-0.929 ± 0.001 -0.930 ± 0.000	$\begin{array}{c} 0.954 \pm 0.001 \\ 0.953 \pm 0.001 \end{array}$	0.996 ± 0.000 0.998 ± 0.000		
			C	IFAR100					
Dropout rate	0.1	MC Dropout Probes + MC Dropout	$\begin{array}{c} 0.044 \pm 0.002 \\ \textbf{0.033} \pm \textbf{0.002} \end{array}$	$\begin{array}{c} 0.754 \pm 0.002 \\ \textbf{0.751} \pm \textbf{0.002} \end{array}$	-0.717 ± 0.002 -0.719 ± 0.002	$\begin{array}{c} 0.802 \pm 0.002 \\ 0.802 \pm 0.002 \end{array}$	$\begin{array}{c} 0.985 \pm 0.000 \\ \textbf{0.988} \pm \textbf{0.000} \end{array}$		
Dropout rate	0.2	MC Dropout Probes + MC Dropout	$\begin{array}{c} 0.035 \pm 0.003 \\ \textbf{0.016} \pm \textbf{0.002} \end{array}$	$\begin{array}{c} 0.709 \pm 0.005 \\ \textbf{0.716} \pm \textbf{0.004} \end{array}$	$\begin{array}{c} -0.718 \pm 0.002 \\ -0.716 \pm 0.002 \end{array}$	$\begin{array}{c} 0.798 \pm 0.002 \\ \textbf{0.794} \pm \textbf{0.002} \end{array}$	$\begin{array}{c} 0.988 \pm 0.000 \\ \textbf{0.991} \pm \textbf{0.001} \end{array}$		
				SVHN					
Dropout rate	0.1	MC Dropout Probes + MC Dropout	$\begin{array}{c} \textbf{0.003} \pm \textbf{0.001} \\ 0.009 \pm 0.001 \end{array}$	$\begin{array}{c} 0.166 \pm 0.010 \\ 0.165 \pm 0.010 \end{array}$	-0.929 ± 0.005 -0.930 ± 0.005	$\begin{array}{c} 0.953 \pm 0.004 \\ 0.954 \pm 0.004 \end{array}$	$\begin{array}{c} 0.996 \pm 0.000 \\ 0.996 \pm 0.000 \end{array}$		
Dropout rate	0.2	MC Dropout Probes + MC Dropout	$\begin{array}{c} 0.016 \pm 0.004 \\ 0.014 \pm 0.002 \end{array}$	0.168 ± 0.007 0.166 ± 0.008	-0.929 ± 0.004 -0.930 ± 0.004	$\begin{array}{c} 0.955 \pm 0.003 \\ 0.955 \pm 0.003 \end{array}$	$\begin{array}{c} 0.997 \pm 0.000 \\ 0.997 \pm 0.000 \end{array}$		
			FASH	ION MNIST					
Dropout rate	0.1	MC Dropout Probes + MC Dropout	$\begin{array}{c} 0.038 \pm 0.005 \\ \textbf{0.011} \pm \textbf{0.001} \end{array}$	$\begin{array}{c} 0.295 \pm 0.033 \\ \textbf{0.237} \pm \textbf{0.015} \end{array}$	-0.857 ± 0.015 -0.881 ± 0.006	$\begin{array}{c} 0.905 \pm 0.010 \\ \textbf{0.918} \pm \textbf{0.004} \end{array}$	$\begin{array}{c} 0.992 \pm 0.001 \\ \textbf{0.996} \pm \textbf{0.001} \end{array}$		
Dropout rate	0.2	MC Dropout Probes + MC Dropout	$\begin{array}{c} 0.018 \pm 0.003 \\ \textbf{0.009} \pm \textbf{0.001} \end{array}$	$\begin{array}{c} 0.233 \pm 0.012 \\ \textbf{0.212} \pm \textbf{0.006} \end{array}$	-0.879 ± 0.005 -0.890 ± 0.003	$\begin{array}{c} 0.917 \pm 0.004 \\ \textbf{0.923} \pm \textbf{0.002} \end{array}$	$\begin{array}{c} 0.996 \pm 0.001 \\ 0.996 \pm 0.000 \end{array}$		

Table 7: Empirical evaluation on Wide-ResNet-28-10 with Monte Carlo dropout. Probe scaling can complement other techniques to improve calibration further.



Figure 3: The relative contribution of each probe is plotted as a function of its relative depth. In classical VGG architectures with no skip connections, probes assigned to the early layers play a significant role in estimating the model's uncertainty. This is not the case in ResNet architectures with skip connections, where only the probes assigned to the upper layers are used for calibration.

5 CONCLUSION

In this paper, we develop "probe scaling", a simple, effective post-hoc technique to improve the calibration of deep neural networks. We empirically demonstrate that using the representations learned by the intermediate layers of a neural network can be utilized to estimate its uncertainty. Besides the theoretical argument for the advantage of probe scaling over temperature scaling, we conduct an empirical evaluation against several post-hoc calibration methods across several architectures and show that probe scaling outperforms other post-hoc methods in standard calibration metrics, such as the Brier score and the negative log-likelihood with statistically significant evidence. In addition, probe scaling is the only post-hoc method that exhibits a consistent improvement over the baseline in both accuracy and AUC. Finally, we show that probe scaling can complement other techniques, such as Monte Carlo dropout, to improve calibration further.

REFERENCES

Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. URL https://www.tensorflow.org/. Software available from tensorflow.org.

- Ibrahim Alabdulmohsin and Mario Lucic. A near-optimal algorithm for debiasing trained machine learning models. *arXiv preprint arXiv:2106.12887*, 2021.
- Devansh Arpit, Stanisław Jastrzkebski, Nicolas Ballas, David Krueger, Emmanuel Bengio, Maxinder S Kanwal, Tegan Maharaj, Asja Fischer, Aaron Courville, Yoshua Bengio, and Simon Lacoste-Julien. A closer look at memorization in deep networks. In *ICML*, 2017.
- Robert JN Baldock, Hartmut Maennel, and Behnam Neyshabur. Deep learning through the lens of example difficulty. *arXiv preprint arXiv:2106.09647*, 2021.
- J Eric Bickel. Some comparisons among quadratic, spherical, and logarithmic scoring rules. *Decision Analysis*, 4(2):49–65, 2007.
- Glenn W Brier et al. Verification of forecasts expressed in terms of probability. *Monthly weather review*, 78(1):1–3, 1950.
- L Elisa Celis, Lingxiao Huang, Vijay Keswani, and Nisheeth K Vishnoi. Classification with fairness constraints: A meta-algorithm with provable guarantees. In *FAccT*, 2019.
- François Chollet et al. Keras. https://keras.io, 2015.
- Gilad Cohen, Guillermo Sapiro, and Raja Giryes. DNN or k-NN: That is the generalize vs. memorize question. *arXiv:1805.06822*, 2018.
- Sam Corbett-Davies, Emma Pierson, Avi Feller, Sharad Goel, and Aziz Huq. Algorithmic decision making and the cost of fairness. In *SIGKDD*, 2017.
- Janez Demšar. Statistical comparisons of classifiers over multiple data sets. *Journal of Machine Learning Research*, 7:1–30, 2006.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Conference on Computer Vision and Pattern Recognition*, 2009.
- Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *ICML*. PMLR, 2016.
- Yonatan Geifman and Ran El-Yaniv. Selective classification for deep neural networks. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *NeurIPS*, 2017.
- Tilmann Gneiting and Adrian E Raftery. Strictly proper scoring rules, prediction, and estimation. *Journal of the American statistical Association*, 102(477):359–378, 2007.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *ICML*. PMLR, 2017.
- Kartik Gupta, Amir Rahimi, Thalaiyasingam Ajanthan, Thomas Mensink, Cristian Sminchisescu, and Richard Hartley. Calibration of neural networks using splines. In *ICLR*, 2021.
- Marton Havasi, Rodolphe Jenatton, Stanislav Fort, Jeremiah Zhe Liu, Jasper Snoek, Balaji Lakshminarayanan, Andrew M Dai, and Dustin Tran. Training independent subnetworks for robust prediction. In *ICLR*, 2021.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *ECCV*, 2016.
- Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. In *International Conference on Learning Representations*, 2019. URL https://openreview.net/forum?id=HJz6tiCqYm.

- Benjamin Kompa, Jasper Snoek, and Andrew L Beam. Second opinion needed: communicating uncertainty in medical machine learning. *NPJ Digital Medicine*, 4(1):1–6, 2021.
- Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, 2009.
- 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. In *NeurIPS*, 2019.
- Aviral Kumar and Sunita Sarawagi. Calibration of encoder decoder models for neural machine translation. arXiv preprint arXiv:1903.00802, 2019.
- Aviral Kumar, Sunita Sarawagi, and Ujjwal Jain. Trainable calibration measures for neural networks from kernel mean embeddings. In *ICML*, pp. 2805–2814. PMLR, 2018.
- Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. *arXiv preprint arXiv:1612.01474*, 2016.
- Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *NeurIPS*, 2017.
- Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv* preprint arXiv:1608.03983, 2016.
- Wesley J Maddox, Pavel Izmailov, Timur Garipov, Dmitry P Vetrov, and Andrew Gordon Wilson. A simple baseline for bayesian uncertainty in deep learning. In *NeurIPS*, 2019.
- Hartmut Maennel, Ibrahim Alabdulmohsin, Ilya Tolstikhin, Robert JN Baldock, Olivier Bousquet, Sylvain Gelly, and Daniel Keysers. What do neural networks learn when trained with random labels? In *NeurIPS*, 2020.
- Aditya Krishna Menon and Robert C Williamson. The cost of fairness in binary classification. In *FAccT*, 2018.
- Zachary Nado, Neil Band, Mark Collier, Josip Djolonga, Michael Dusenberry, Sebastian Farquhar, Angelos Filos, Marton Havasi, Rodolphe Jenatton, Ghassen Jerfel, Jeremiah Liu, Zelda Mariet, Jeremy Nixon, Shreyas Padhy, Jie Ren, Tim Rudner, Yeming Wen, Florian Wenzel, Kevin Murphy, D. Sculley, Balaji Lakshminarayanan, Jasper Snoek, Yarin Gal, and Dustin Tran. Uncertainty Baselines: Benchmarks for uncertainty & robustness in deep learning. arXiv preprint arXiv:2106.04015, 2021.
- Mahdi Pakdaman Naeini, Gregory F. Cooper, and Milos Hauskrecht. Obtaining well calibrated probabilities using bayesian binning. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, AAAI'15, pp. 2901–2907. AAAI Press, 2015. ISBN 0262511290.
- Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. 2011.
- Jeremy Nixon, Mike Dusenberry, Ghassen Jerfel, Timothy Nguyen, Jeremiah Liu, Linchuan Zhang, and Dustin Tran. Measuring calibration in deep learning. arXiv preprint arXiv:1904.01685, 2019.
- Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, D. Sculley, Sebastian Nowozin, Joshua Dillon, Balaji Lakshminarayanan, and Jasper Snoek. Can you trust your model's uncertainty? evaluating predictive uncertainty under dataset shift. In *NeurIPS*, 2019.
- John Platt et al. Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. *Advances in large margin classifiers*, 10(3):61–74, 1999.
- Maithra Raghu, Chiyuan Zhang, Jon Kleinberg, and Samy Bengio. Transfusion: Understanding transfer learning for medical imaging. In *NeurIPS*, 2019.
- Shai Shalev-Shwartz and Shai Ben-David. Understanding machine learning: From theory to algorithms. Cambridge university press, 2014.

- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.
- Sunil Thulasidasan, Gopinath Chennupati, Jeff A Bilmes, Tanmoy Bhattacharya, and Sarah Michalak. On mixup training: Improved calibration and predictive uncertainty for deep neural networks. In *NeurIPS*, 2019.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *CoRR*, abs/1708.07747, 2017. URL http://arxiv.org/ abs/1708.07747.
- Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural networks? In *NeurIPS*, 2014.
- Bianca Zadrozny and Charles Elkan. Transforming classifier scores into accurate multiclass probability estimates. In SIGKDD, 2002.
- Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.

A PROOF OF PROPOSITION 1

The proof of this proposition uses Corollary 13.8 in Shalev-Shwartz & Ben-David (2014). Let $L_S(w)$ be the empirical loss w.r.t. f and write $L_D(w)$ for the population loss of w. If f is a convex function that is ρ -Lipschitz continuous, then the solution to the regularized empirical risk minimization problem:

$$\hat{\mathbf{w}} = \arg\min_{w} \left\{ \lambda ||w||^2 + L_S(w) \right\},\,$$

satisfies the oracle inequality:

$$\forall \mathbf{w}: \quad \mathbb{E}_{S} \left[L_{\mathcal{D}}(\hat{\mathbf{w}}] \leq L_{\mathcal{D}}(\mathbf{w}) + \lambda ||\mathbf{w}||^{2} + \frac{2\rho^{2}}{\lambda N} \right]$$

Choosing $\mathbf{w} = \mathbf{w}^{\star} \doteq \arg \min_{w} \{ L_{\mathcal{D}}(w) \}$:

$$\mathbb{E}_{S}\left[L_{\mathcal{D}}(\hat{\mathbf{w}}] \leq L_{\mathcal{D}}(\mathbf{w}^{\star}) + \lambda ||\mathbf{w}^{\star}||^{2} + \frac{2\rho^{2}}{\lambda N}\right]$$

Using the definition of \mathbf{w}^* and applying Markov's inequality, we have for any \mathbf{w} :

$$p_{S}\left\{L_{\mathcal{D}}(\hat{\mathbf{w}}) - L_{\mathcal{D}}(\mathbf{w}) \ge \epsilon\right\} \le p_{S}\left\{L_{\mathcal{D}}(\hat{\mathbf{w}}) - L_{\mathcal{D}}(\mathbf{w}^{\star}) \ge \epsilon\right\}$$
$$\le \frac{1}{\epsilon} \left(\lambda ||\mathbf{w}^{\star}||^{2} + \frac{2\rho^{2}}{\lambda N}\right)$$

Setting $\hat{\mathbf{w}}$ for the solution learned by probe scaling and \mathbf{w} for temperature scaling (which is valid because temperature scaling belongs to the search space of probe scaling), we obtain the statement of the proposition.