

# ESD: EXPECTED SQUARED DIFFERENCE AS A TUNING-FREE TRAINABLE CALIBRATION MEASURE

Hee Suk Yoon<sup>1\*</sup> Joshua Tian Jin Tee<sup>1\*</sup> Eunseop Yoon<sup>1</sup> Sunjae Yoon<sup>1</sup>  
 Gwangsu Kim<sup>1</sup> Yingzhen Li<sup>2</sup> Chang D. Yoo<sup>1†</sup>  
<sup>1</sup>Korea Advanced Institute of Science and Technology (KAIST) <sup>2</sup>Imperial College London  
 {hskyo, joshuateetj, esyo, sunjae.yoon}@kaist.ac.kr  
 s88012@gmail.com yingzhen.li@imperial.ac.uk cd\_yoo@kaist.ac.kr

## ABSTRACT

Studies have shown that modern neural networks tend to be poorly calibrated due to over-confident predictions. Traditionally, post-processing methods have been used to calibrate the model after training. In recent years, various trainable calibration measures have been proposed to incorporate them directly into the training process. However, these methods all incorporate internal hyperparameters, and the performance of these calibration objectives relies on tuning these hyperparameters, incurring more computational costs as the size of neural networks and datasets become larger. As such, we present Expected Squared Difference (ESD), a tuning-free (i.e., hyperparameter-free) trainable calibration objective loss, where we view the calibration error from the perspective of the squared difference between the two expectations. With extensive experiments on several architectures (CNNs, Transformers) and datasets, we demonstrate that (1) incorporating ESD into the training improves model calibration in various batch size settings without the need for internal hyperparameter tuning, (2) ESD yields the best-calibrated results compared with previous approaches, and (3) ESD drastically improves the computational costs required for calibration during training due to the absence of internal hyperparameter. The code is publicly accessible at <https://github.com/hee-suk-yoon/ESD>.

## 1 INTRODUCTION

The calibration of a neural network measures the extent to which its predictions align with the true probability distribution. Possessing this property becomes especially important in real-world applications, such as identification (Kim & Yoo, 2017; Yoon et al., 2022), autonomous driving (Bojarski et al., 2016; Ko et al., 2017), and medical diagnosis (Kocbek et al., 2020; Pham et al., 2022), where uncertainty-based decisions of the neural network are crucial to guarantee the safety of the users. However, despite the success of modern neural networks in accurate classification, they are shown to be poorly calibrated due to the tendency of the network to make predictions with high confidence regardless of the input. (i.e., over-confident predictions) (Guo et al., 2017).

Traditionally, post-processing methods have been used, such as temperature scaling and vector scaling (Guo et al., 2017), to calibrate the model using the validation set after the training by adjusting the logits before the final softmax layer. Various trainable calibration objectives have been proposed recently, such as MMCE (Kumar et al., 2018) and SB-ECE (Karandikar et al., 2021), which are added to the loss function as a regularizer to jointly optimize accuracy and calibration during training. A key advantage of calibration during training is that it is possible to cascade post-processing calibration methods after training to achieve even better-calibrated models. Unfortunately, these existing approaches introduce additional hyperparameters in their proposed calibration objectives, and the performance of the calibration objectives is highly sensitive to these design choices. Therefore these hyperparameters need to be tuned carefully on a per model per dataset basis, which greatly reduces their viability for training on large models and datasets.

---

\*Equal contribution

†Corresponding Author

To this end, we propose Expected Squared Difference (ESD), a trainable calibration objective loss that is hyperparameter-free. ESD is inspired by the KS-Error (Gupta et al., 2021), and it views the calibration error from the perspective of the difference between the two expectations. In detail, our contributions can be summarized as follows:

- We propose ESD as a trainable calibration objective loss that can be jointly optimized with the negative log-likelihood loss (NLL) during training. ESD is a binning-free calibration objective loss, and no additional hyperparameters are required. We also provide an unbiased and consistent estimator of the Expected Squared Difference, and show that it can be utilized in small batch train settings.
- With extensive experiments, we demonstrate that across various architectures (CNNs & Transformers) and datasets (in vision & NLP domains), ESD provides the best calibration results when compared to previous approaches. The calibrations of these models are further improved by post-processing methods.
- We show that due to the absence of an internal hyperparameter in ESD that needs to be tuned, it offers a drastic improvement compared to previous calibration objective losses with regard to the total computational cost for training. The discrepancy in computational cost between ESD and tuning-required calibration objective losses becomes larger as the model complexity and dataset size increases.

## 2 RELATED WORK

Calibration of neural networks has gained much attention following the observation from Guo et al. (2017) that modern neural networks are poorly calibrated. One way to achieve better calibration is to design better neural network architectures tailored for uncertainty estimation, e.g., Bayesian Neural Networks (Blundell et al., 2015; Gal & Ghahramani, 2016) and Deep Ensembles (Lakshminarayanan et al., 2017). Besides model design, post-processing calibration strategies have been widely used to calibrate a trained machine learning model using a hold-out validation dataset. Examples include temperature scaling (Guo et al., 2017), which scales the logit output of a classifier with a temperature parameter; Platt scaling (Platt, 1999), which fits a logistic regression model on top of the logits; and Conformal prediction (Vovk et al., 2005; Lei et al., 2018) which uses validation set to estimate the quantiles of a given scoring function. Other post-processing techniques include histogram binning (Zadrozny & Elkan, 2001), isotonic regression (Zadrozny & Elkan, 2002), and Bayesian binning into quantiles (Pakdaman Naeini et al., 2015).

Our work focuses on trainable calibration methods, which train neural networks using a hybrid objective, combining a primary training loss with an auxiliary calibration objective loss. In this regard, one popular objective is Maximum Mean Calibration Error (MMCE) (Kumar et al., 2018), which is a kernel embedding-based measure of calibration that is differentiable and, therefore, suitable as a calibration loss. Moreover, Karandikar et al. (2021) proposes a trainable calibration objective loss, SB-ECE, and S-AvUC, which softens previously defined calibration measures.

## 3 PROBLEM SETUP

### 3.1 CALIBRATION ERROR AND METRIC

Let us first consider an arbitrary neural network  $f: D \rightarrow [0, 1]^C$  with network parameters  $\theta$ , where  $D$  is the input domain and  $C$  is the number of classes in the multiclass classification task. Furthermore, we assume that the training data  $(x_i; y_i)_{i=1}^n$  are sampled i.i.d. from the joint distribution  $P(X; Y)$  (here we use one-hot vector  $y$ ). Here,  $Y = (Y_1; \dots; Y_C)$ , and  $y = (y_1; \dots; y_C)$  are samples from this distribution. We can further define a multivariate random variable  $f(X)$  as the distribution of the outputs of the neural network. Similarly,  $Z = (Z_1; \dots; Z_C)$ , and  $z = (z_1; \dots; z_C)$  are samples from this distribution. We use  $(z_{k,i}; y_{k,i})$  to denote the output confidence and the one-hot vector element associated with the  $k$ -th class of the  $i$ -th training sample. Using this formulation, a neural network is said to be perfectly calibrated for class  $k$  if and only if

$$P(Y_k = 1 | Z_k = z_k) = z_k \tag{1}$$

Intuitively, Eq. (1) requires the model accuracy for class  $k$  to be  $z_k$  on average for the inputs where the neural network produces prediction of class  $k$  with confidence  $z_k$ . In many cases, the research of calibration mainly focuses on the calibration of the max output class that the model predicts. Thus, calibration error is normally reported with respect to the predicted class (i.e., max output class) only. As such,  $k$  will denote the class with maximum output probability. Furthermore, we also write  $\mathbb{1}_{\{x\}}$  as the indicator function which returns one if the Boolean expression is true and zero otherwise.

With these notations, one common measurement of calibration error is the difference between confidence and accuracy which is mathematically represented as (Guo et al., 2017)

$$E_{Z_k} [\mathbb{1}_{\{P(Y_k = 1|Z_k) - Z_k\}}] \quad (2)$$

To estimate this, the Expected Calibration Error (ECE) (Naeini et al., 2015) uses number of bins with disjoint intervals  $B_j = (\frac{j}{B}, \frac{j+1}{B}]$ ,  $j = 0; 1; \dots; B-1$  to compute the calibration error as follows:

$$ECE = \frac{1}{|D|} \sum_{j=0}^{B-1} \sum_{i=1}^N \mathbb{1}_{\{\frac{j}{B} < z_{k;i} \leq \frac{j+1}{B}; y_{k;i} \neq 1\}} \quad (3)$$

### 3.2 CALIBRATION DURING TRAINING

Post-processing method and calibration during training are the two primary approaches for calibrating a neural network. Our focus in this paper is on the latter, where our goal is to train a calibrated yet accurate classifier directly. Note that calibration and predictive accuracy are independent properties of a classifier. In other words, being calibrated does not imply that the classifier has good accuracy and vice versa. Thus, training a classifier to have both high accuracy and good calibration requires jointly optimizing a calibration objective loss alongside the negative log-likelihood (NLL) with a scaling parameter for the secondary objective:

$$\min \text{NLL}(D; \theta) + \lambda \text{CalibrationObjective}(\theta; \gamma) \quad (4)$$

#### 3.2.1 EXISTING TRAINABLE CALIBRATION OBJECTIVES NEED TUNING

Kumar et al. (2018) and Karandikar et al. (2021) suggested that the disjoint bins in ECE can introduce discontinuities which are problematic when using it as a calibration loss in training. Therefore, additional parameters were introduced for the purpose of calibration during training. For example, Kumar et al. (2018) proposed Maximum Mean Calibration Error (MMCE), where it utilizes a Laplacian Kernel instead of the disjoint bins in ECE. Furthermore, Karandikar et al. (2021) proposed soft-binned ECE (SB-ECE) and soft AvUC, which are softened versions of the ECE metric and the AvUC loss (Krishnan & Tickoo, 2020), respectively. All these approaches address the issue of discontinuity which makes the training objective differentiable. However they all introduce additional design choices - MMCE requires a careful selection of the kernel width  $\sigma$ , SB-ECE needs to choose the number of bins  $M$  and a softening parameter  $\tau$ , and S-AvUC requires a user-specified entropy threshold  $\epsilon$  in addition to the softening parameter  $\tau$ . Searching for the optimal hyperparameters can be computationally expensive especially as the size of models and dataset become larger.

#### 3.2.2 CALIBRATION DURING TRAINING SUFFERS FROM OVER-FITTING PROBLEM

NLL loss is known to implicitly train for calibration since it is a proper scoring rule, so models trained with NLL can overfit in terms of calibration error (Mukhoti et al., 2020; Karandikar et al., 2021). Figure 1 provides an example of an accuracy curve and its corresponding ECE curve for a model trained with NLL loss. We see that the model is overfitting in terms of both accuracy and ECE, which causes the gap between train and test ECE to become larger during training. As such, adding a calibration objective to the total loss will not be able to improve model calibration as it mainly helps improve calibration of the model with respect to the training data only.

Kumar et al. (2018) tackle this issue by introducing a weighted version of their calibration objective loss. They consider using larger weights to incorrect prediction samples after observing that the fraction of incorrect to the correct samples on the training data is smaller than that of the validation and test data. Instead of changing the calibration objective function itself, Karandikar et al. (2021) introduced a new training scheme called interleaved training where they split the training data

Figure 1: Accuracy (%) curve (left) and its corresponding ECE (%) curve (right) during training with negative log-likelihood (NLL) loss. It could be seen that since NLL implicitly trains for calibration error, the ECE of the train set approaches zero while the ECE of the test set increases during training.

further and dedicate a small portion of it to optimize the calibration objective loss. We will follow this strategy and introduce the technical details in the later sections.

#### 4 ESD: EXPECTED SQUARED DIFFERENCE

We propose Expected Squared Difference (ESD) as a tuning-free (i.e., hyperparameter-free) calibration objective. Our approach is inspired by the viewpoint of calibration error as a measure of distance between two distributions to obtain a binning-free calibration metric (Gupta et al., 2021). By using this approach, there is no need to employ kernels or softening operations to handle the bins within calibration metrics, such as ECE, in order to make them suitable for training.

In particular, we consider calibration error as the difference between the two expectations. We start with the definition of perfect calibration in Eq. (1):

$$\begin{aligned} P(Y_k = 1 | Z_k = z_k) &= z_k \quad z_k \in [0, 1] \\ P(Y_k = 1; Z_k = z_k) &= z_k P(Z_k = z_k) \quad z_k \in [0, 1] \quad (\text{by Bayes rule}) \end{aligned} \quad (5)$$

Now considering the accumulation of the terms on both sides for arbitrary confidence  $z_k \in [0, 1]$ , the perfect calibration for an arbitrary class can now be written as:

$$\begin{aligned} \int_0^1 P(Y_k = 1; Z_k = z_k) dz_k &= \int_0^1 z_k P(Z_k = z_k) dz_k \\ E_{Z_k; Y_k} [I(Z_k = z_k; Y_k = 1)] &= E_{Z_k; Y_k} [Z_k I(Z_k = z_k)]; \quad z_k \in [0, 1]; \end{aligned} \quad (6)$$

This allows us to write the difference between the two expectations as

$$\begin{aligned} d_k(z_k) &= \int_0^1 P(Y_k = 1; Z_k = z_k) - z_k P(Z_k = z_k) dz_k \\ &= \int E_{Z_k; Y_k} [I(Z_k = z_k; Y_k = 1) - I(Z_k = z_k)]; \end{aligned} \quad (7)$$

and  $d_k(z_k) = 0, z_k \in [0, 1]$  if and only if the model is perfectly calibrated for class  $k$ . Since this has to hold  $z_k \in [0, 1]$ , we propose ESD as the expected squared difference between the two expectations:

$$E_{Z_k^0} [d_k(Z_k^0)^2] = E_{Z_k^0} [E_{Z_k; Y_k}^2 [I(Z_k = Z_k^0; Y_k = 1) - I(Z_k = Z_k^0)]]; \quad (8)$$

Due to the close relationship between  $d_k(z_k)$  and calibration of a neural network, the difference between an uncalibrated and a calibrated neural network can be clearly observed by using  $d_k(Z_k^0)^2$  as visually shown in Appendix A. Since  $ESD = 0$  iff. the model is perfectly calibrated as shown in the following theorem, this metric is a good measure of calibration.

**Theorem 1.**  $E_{Z_k^0} [E_{Z_k; Y_k}^2 [I(Z_k = Z_k^0; Y_k = 1) - I(Z_k = Z_k^0)]] = 0$  iff. the model is perfectly calibrated.

**Proof.** Since  $d_k(Z_k^0)^2 = E_{Z_k; Y_k}^2 [I(Z_k = Z_k^0; Y_k = 1) - I(Z_k = Z_k^0)]$  is a non-negative random variable induced by  $Z_k^0$ ,  $E_{Z_k^0} [d_k(Z_k^0)^2] = 0$  iff.  $P(d_k(Z_k^0) = 0) = 1$ : Furthermore,

$$P(d_k(Z_k^0) = 0) = 1 \iff d_k(z_k) = 0 \quad \forall z_k \in \text{supp}(Z_k^0);$$

Thus,

$$E_{Z_k^0}[d_k(Z_k^0)^2] = 0 \text{ iff. } d_k(\cdot) = 0 \quad \forall z_k \in \mathcal{Z}_k$$

From lemma 1.1 below, we have  $d_k(\cdot) = 0 \quad \forall z_k \in \mathcal{Z}_k$  iff.  $d_k(\cdot) = 0 \quad \forall z_k \in \mathcal{Z}_k$  s.t.  $P(Z_k^0 = z_k) > 0$ . Consequently  $E_{Z_k^0}[E_{Z_k; Y_k}^2[I(Z_k = Z_k^0)(I(Y_k = 1) - Z_k)]] = 0$  iff. the model is perfectly calibrated.  $\square$

**Lemma 1.1** Let  $\mathcal{Z}_k$  be the support set of random variable  $Z_k$ , then  $d_k(\cdot) = 0 \quad \forall z_k \in \mathcal{Z}_k$  iff.  $d_k(\cdot) = 0 \quad \forall z_k \in \mathcal{Z}_k$ .

*Proof.* The backward direction result is straight-forward, so we only prove for the forward direction: If  $d_k(\cdot) = 0 \quad \forall z_k \in \mathcal{Z}_k$ , then  $d_k(\cdot) = 0 \quad \forall z_k \in \mathcal{Z}_k$ .

For arbitrary  $z_k^0 \in \mathcal{Z}_k$ , let  $z_k^* = \arg \min_{z_k \in \mathcal{Z}_k} |z_k - z_k^0|$  and  $\epsilon > 0$ . We then have,

$$\begin{aligned} d_k(z_k^0) &= \int_{\mathcal{Z}_k} P(Y_k = 1; Z_k = z_k) - z_k P(Z_k = z_k) dz_k \\ &= \int_{\mathcal{Z}_k} P(Y_k = 1 | Z_k = z_k) P(Z_k = z_k) - z_k P(Z_k = z_k) dz_k \\ &= \int_{\mathcal{Z}_k} P(Y_k = 1 | Z_k = z_k) P(Z_k = z_k) - z_k P(Z_k = z_k) dz_k \\ &\quad + \int_{\mathcal{Z}_k} P(Y_k = 1 | Z_k = z_k) P(Z_k = z_k) - z_k P(Z_k = z_k) dz_k \\ &= \int_{\mathcal{Z}_k} P(Y_k = 1 | Z_k = z_k) P(Z_k = z_k) - z_k P(Z_k = z_k) dz_k \\ &= 0: \end{aligned}$$

$\square$

#### 4.1 AN ESTIMATOR FORESD

In this section, we use  $(z_{k,i}; y_{k,i})$  to denote the output confidence and the one-hot vector element associated with the  $k$ -th class of the  $i$ -th training sample respectively. As the expectations in the true Expected Squared Difference (Eq. (8)) are intractable, we propose a Monte Carlo estimator for it which is unbiased. A common approach is to use a naive Monte Carlo sampling with respect to both the inner and outer expectations to give the following:

$$E_{Z_k^0}[E_{Z_k; Y_k}^2[I(Z_k = Z_k^0)(I(Y_k = 1) - Z_k)]] = \frac{1}{N} \sum_{i=1}^N g_i^2; \quad (9)$$

$$\text{where } g_i = \frac{1}{N-1} \sum_{\substack{j=1 \\ j \neq i}}^N g_{ij} \text{ and } g_{ij} = I(z_{k,j} = z_{k,i}) [I(y_j = 1) - z_{k,j}];$$

However, Eq. (9) results in a biased estimator that is an upper bound of the true Expected Squared Difference. To account for the bias, we propose the unbiased and consistent estimator (proof in Appendix B) of the true Expected Squared Difference:

$$\text{ESD} = \frac{1}{N} \sum_{i=1}^N g_i^2 - \frac{S_{g_i}^2}{N-1} \quad \text{where } S_{g_i}^2 = \frac{1}{N-2} \sum_{\substack{j=1 \\ j \neq i}}^N (g_{ij} - g_i)^2; \quad (10)$$

<sup>1</sup>To avoid confusion, ESD from this point onward will refer to the estimator instead of its expectation form.

## 4.2 INTERLEAVED TRAINING

Negative log-likelihood (NLL) has been shown to greatly overfit to ECE of the data it is trained on. Thus, training for calibration using the same data for the NLL has limited effect on reducing the calibration error of the model. Karandikar et al. (2021) proposed interleaved training where they split the train set into two subsets - one is used to optimize the NLL and the other is used to optimize the calibration objective. Following this framework,  $\mathcal{D}_{\text{train}}$  denote the entire train set. We separate the train set into two subsets  $\mathcal{D}_{\text{train}}^0$  and  $\mathcal{D}_{\text{cal}}^0$ . The joint training of ESD with NLL becomes,

$$\min \text{NLL}(\mathcal{D}_{\text{train}}^0; \theta) + \text{ESD}(\mathcal{D}_{\text{cal}}^0; \theta); \quad (11)$$

With this training scheme, NLL is optimized on  $\mathcal{D}_{\text{train}}^0$  and ESD is optimized on  $\mathcal{D}_{\text{cal}}^0$ . This way, Karandikar et al. (2021) showed we can avoid minimizing ECE that is already overfitted to the train set.

## 5 EXPERIMENTAL SETTING

### 5.1 DATASETS AND MODELS

**Image Classification** For image classification tasks we use the following datasets:

- MNIST (Deng, 2012): 54,000/6,000/10,000 images for train, validation, and test split was used. We resized the images to (32x32) before inputting to the network.
- CIFAR10 & CIFAR100 (Krizhevsky et al., 2010): 45,000/5,000/10,000 images for train, validation, and test split was used. Used random cropping of 32 with padding of 4. Normalized each RGB channel with mean of 0.5 and standard deviation of 0.5.
- ImageNet100 (Deng et al., 2009): A subset dataset from ImageNet Large Scale Visual Recognition Challenge 2012 with 100 classes. Since the labels of test sets are unavailable, we use the official validation set as the test set, and we dedicate 10% of the training data as the validation set for the experiments (i.e., 117,000/13,000/5,000 split for train/val/test set).

**Natural Language Inference (NLI)** NLI is a task in Natural Language Processing where it involves classifying the inference relation (entailment, contradiction, or neutral) between two texts (MacCartney & Manning, 2008). For NLI tasks we use the following datasets:

- SNLI (Bowman et al., 2015): SNLI corpus is a collection of human-written English sentence pairs manually labeled for balanced classification with the labels entailment, contradiction, and neutral. The data consists of 550,152/10,000/10,000 sentence pairs for train/val/test set respectively. The max length of the input is set to 158.
- ANLI (Nie et al., 2020): ANLI dataset is a large-scale NLI dataset, collected via an adversarial human-and-model-in-the-loop procedure. The data consists of 162,865/3,200/3,200 sentence pairs for train/val/test set respectively. The max length of the input is set to 128.

For the Image Classification datasets, we used Convolutional Neural Networks (CNNs). Specifically, we used LeNet5 (Lecun et al., 1998), ResNet50, ResNet34, and ResNet18 (He et al., 2016), for MNIST, CIFAR10, CIFAR100, and ImageNet100, respectively. For the NLI datasets, we used pretrained transformer based Pre-trained Language Models (PLMs). Specifically, we used Bert-base (Devlin et al., 2019) and Roberta-base (Liu et al.), for SNLI and ANLI, respectively.

### 5.2 EXPERIMENTAL SETUP

We compare our Expected Squared Difference (ESD) to previously proposed trainable calibration objectives, MMCE and SB-ECE, on the datasets and models previously mentioned. For fair comparison, interleaved training has been used for all three calibration objectives. For MMCE, in which Kumar et al. (2018) proposed an unweighted and weighted version of the objective, we use the former to set the method to account for the overfitting problem mentioned in section 3.2.2 consistent to interleaved training. For the interleaved training settings, we held out 10% of the train set to the calibration set. The regularizer hyperparameter for weighting the calibration measure with respect to NLL is chosen via a grid search. For measuring calibration error, we use ECE with 20 equally sized

<sup>2</sup>For  $\lambda$ , we search for [0.2, 0.4, 0.6, 0.8, 1.0, 2.0, 3.0, ... 10.0] (Appendix E).

Dataset (Model)	Loss Fn.	Acc.		ECE		ECE after TS		Acc. after VS		ECE after VS	
MNIST (LeNet5)	NLL (baseline)	98.8	0.034	0.91	0.080	0.31	0.044	98.7	0.058	0.43	0.123
	+MMCE	98.4	0.320	0.36	0.029	0.33	0.068	98.4	0.250	0.32	0.055
	+SB-ECE	97.9	0.177	0.41	0.067	0.45	0.103	97.9	0.073	0.39	0.088
	+ESD (ours)	98.6	0.204	0.30	0.035	0.29	0.030	98.6	0.167	0.28	0.071
CIFAR10 (Resnet50)	NLL (baseline)	92.9	0.159	5.49	0.105	1.89	0.105	92.3	0.625	2.96	1.596
	+MMCE	91.5	0.340	4.92	0.292	1.92	0.274	91.3	0.429	2.52	0.586
	+SB-ECE	91.6	0.288	4.86	0.319	1.63	0.300	91.5	0.343	1.88	0.393
	+ESD (ours)	92.1	0.141	3.08	0.692	1.60	0.089	92.1	0.314	1.61	0.287
CIFAR100 (Resnet34)	NLL (baseline)	68.4	0.491	23.8	0.403	5.74	0.306	67.3	0.551	9.78	0.640
	+MMCE	66.5	0.644	13.9	0.545	4.91	0.457	66.3	0.773	4.59	1.575
	+SB-ECE	67.3	0.367	14.7	1.370	4.94	0.240	66.4	0.369	4.57	1.253
	+ESD (ours)	67.4	0.356	13.6	0.950	4.85	0.390	67.1	0.496	4.28	1.396
ImageNet100 (Resnet18)	NLL (baseline)	75.8	0.397	10.3	0.686	2.51	0.378	75.9	0.595	2.86	0.459
	+MMCE	74.3	0.248	3.83	0.644	1.94	0.262	74.5	0.444	2.29	0.227
	+SB-ECE	74.4	0.596	4.28	0.318	2.06	0.260	74.7	0.472	2.13	0.198
	+ESD (ours)	74.6	0.320	1.80	0.262	1.72	0.350	74.8	0.436	1.87	0.266
SNLI (Bert-base)	NLL (baseline)	90.2	0.434	4.20	0.420	1.04	0.112	90.1	0.397	0.96	0.099
	+MMCE	89.4	0.491	1.11	0.181	1.01	0.178	89.5	0.588	1.08	0.170
	+SB-ECE	89.1	0.668	1.82	0.301	0.99	0.180	89.2	0.765	0.90	0.109
	+ESD (ours)	89.3	0.536	0.98	0.165	0.73	0.192	89.4	0.583	0.61	0.126
ANLI (Roberta-base)	NLL (baseline)	49.4	0.323	35.9	0.505	4.16	0.445	48.4	0.492	5.10	0.657
	+MMCE	49.0	0.471	31.2	0.850	3.71	0.175	47.7	0.429	4.79	0.693
	+SB-ECE	48.5	0.481	33.9	1.378	3.98	0.733	47.3	0.281	5.10	0.119
	+ESD (ours)	48.0	0.451	28.8	0.543	3.49	0.373	47.1	0.429	4.42	1.010

Table 1: Average accuracy (%) and ECE (%) (with std. across 5 trials) for baseline, MMCE, SB-ECE, ESD after training and after post-processing with temperature scaling (TS) or vector scaling (VS).

bins. For the image classification tasks we use AdamW (Loshchilov & Hutter, 2019) optimizer with  $10^{-3}$  learning rate and  $10^{-2}$  weight decay for 250 epochs, except for ImageNet100, in which case we used  $10^{-4}$  weight decay for 90 epochs. For the NLI tasks, we use AdamW optimizer with  $10^{-3}$  learning rate and  $10^{-2}$  weight decay for 15 epochs. For both tasks, we use a batch size of 512.

The internal hyperparameters within MMCE (and SB-ECE) ( $\alpha, \beta, \gamma, \delta, \epsilon, \eta, \theta, \tau$ ) were sequentially optimized<sup>3</sup> following the search for optimal. Following the original experimental setting of SB-ECE by Karandikar et al. (2021), we fix the hyperparameter  $\tau$  to 15. Similar to the model selection criterion utilized by Karandikar et al. (2021), we look at the accuracy and the ECE of all possible hyperparameter configurations in the grid. We choose the lowest ECE while giving up less than 1.5% accuracy relative to baseline accuracy on the validation set. All experiments were done using NVIDIA Quadro RTX 8000 and NVIDIA RTX A6000.

## 6 EXPERIMENTAL RESULT

In Table 1, we report the accuracy and ECE for the models before and after post-processing (i.e., temperature scaling and vector scaling) of various datasets and models. Compared to the baseline or other trainable calibration objective loss (MMCE and SB-ECE), jointly training with ESD as the secondary loss consistently results in a better-calibrated network for all the datasets and models with around 1% degradation of accuracy. Moreover, we observe that applying post-processing (i.e., temperature scaling (TS) and vector scaling (VS)) after training with a calibration objective loss generally results in better-calibrated models over baseline post-processing, in which case ESD still outperforms other methods. Comparing the calibration outcomes in temperature scaling with vector scaling for models trained with ESD, vector scaling worked comparable if not better as a post-processing method, except for ANLI, with minimal impact on the accuracy for all datasets. Additionally, Ovadia et al. (2019) demonstrated that postprocessing methods, such as temperature

<sup>3</sup>For  $\alpha$ , we search [0.2, 0.4, 0.6, 0.8]. For  $\beta$ , we search [0.0001, 0.001, 0.01, 0.1].

Figure 2: ECE performance curve of MMCE (left) and SB-ECE (right) with respect to their varying internal hyperparameters on MNIST, CIFAR10, SNLI datasets.

Figure 3: Computational cost of single-run training (left) and total cost considering hyperparameter tuning (right). The x-axis in both cases are in the order of increasing model complexity.

scaling, are not effective in achieving good calibration under distribution shifts. Building on this research, Karandikar et al. (2021) have provided evidence that incorporating calibration during model training can enhance a model's ability to maintain calibration under distribution shifts. This finding has been replicated in Appendix F.

We show in Figure 2 that across different models and datasets the calibration performance of MMCE and SB-ECE is sensitive to the internal hyperparameter. This shows the importance of hyperparameter tuning in these methods. On the other hand, ESD does not only outperforms MMCE and SB-ECE in terms of ECE but does not have an internal hyperparameter to tune for.

## 7 ABLATION STUDY

### 7.1 COMPUTATIONAL COST OF HYPERPARAMETER SEARCH

We investigate the computational cost required to train a model with ESD compared to MMCE and SB-ECE. As depicted in Figure 3, the computational cost required to train a model for a single run remains nearly identical across different models and datasets. However, considering the need for tuning additional hyperparameters within MMCE and SB-ECE, the discrepancy in total computational cost between ESD and tuning-required calibration objective losses becomes more prominent as the model complexity and dataset size increases.

### 7.2 BATCH SIZE EXPERIMENTS

Table 2 shows that ESD is robust to varying batch sizes which could be due to the fact that it is an unbiased and consistent estimator. Regardless of the batch size, being unbiased guarantees that the average of the gradient vectors of ESD equals the gradient of the true estimate (Appendix D).



Dataset (Model)	Loss Fn.	Acc.	ECE	ECE after TS	Acc. after VS	ECE after VS
CIFAR10 (Resnet50)	NLL (b512)	92.9 0.159	5.49 0.105	1.89 0.105	92.3 0.625	2.96 1.596
	+ESD (b512)	92.1 0.141	3.08 0.692	1.60 0.089	92.1 0.314	1.61 0.287
	+ESD (b256)	91.5 0.137	3.15 0.587	1.61 0.073	91.9 0.573	1.60 0.082
	+ESD (b128)	91.5 0.519	2.93 0.539	1.70 0.600	91.7 0.246	2.20 0.584
	+ESD (b64)	91.3 0.753	3.91 0.847	1.87 0.132	91.5 0.257	2.65 1.48
SNLI (Bert-base)	NLL (b512)	90.2 0.434	4.20 0.420	1.04 0.112	90.1 0.397	0.96 0.099
	+ESD (b512)	89.3 0.536	0.98 0.165	0.73 0.192	89.4 0.583	0.61 0.126
	+ESD (b256)	89.9 0.314	1.11 0.236	0.76 0.241	88.9 0.146	0.62 0.147
	+ESD (b128)	89.8 0.435	0.99 0.244	0.68 0.190	89.6 0.490	0.58 0.187
	+ESD (b64)	89.1 0.660	1.06 0.285	0.83 0.260	88.8 0.672	0.81 0.153

Table 2: Average accuracy and ECE (with std. across 5 trials) for ESD after training with batch sizes 64, 128, 256, and 512.

### 7.3 ARE INDICATOR FUNCTIONS TRAINABLE ?

Recent papers, Kumar et al. (2018) and Karandikar et al. (2021), have suggested that ECE is not suitable for training as a result of its high discontinuity due to binning, which can be seen as a form of an indicator function. However, the results from our method suggest that our measure was still able to train well despite the existence of indicator functions. In addition, previous measures also contain indicator functions in the form of  $\arg\max$  function that introduces discontinuities but remains to be trainable. As such, this brings to rise the question of whether calibration measures with indicator functions can be used for training. To investigate this, we ran ECE as an auxiliary calibration loss on different batch sizes and observed its performance on CIFAR 100 (Table 3).

We found that ECE, contrary to previous belief, is trainable under large batch sizes and not trainable under small batch sizes while ours maintains good performance regardless of batch size. The poor performance of ECE under small batch size setting could potentially be attributed to the high bias present in such cases. From this, it seems to suggest that indicator functions do not seem to inhibit the training for calibration.

Loss Fn.	Acc.	ECE
NLL (b512)	68.6 0.034	22.4 0.105
+ECE (b512)	68.2 0.025	14.6 0.252
+ECE (b256)	67.9 0.057	18.1 0.552
+ECE (b128)	68.1 0.084	21.0 0.552

Table 3: Average accuracy and ECE (with std. across 5 trials) for CIFAR100 after training with ECE on batch sizes 128, 256, and 512.

## 8 CONCLUSIONS

Motivated by the need for a tuning-free trainable calibration objective, we proposed Expected Squared Difference (ESD), which does not contain any internal hyperparameters. With extensive comparison with existing methods for calibration during training across various architectures (CNNs & Transformers) and datasets (in vision & NLP domains), we demonstrate that training with ESD provides the best-calibrated models compared to other methods with minor degradation in accuracy over baseline. Furthermore, the calibration of these models is further improved after post-processing. In addition, we demonstrate that ESD can be utilized in small batch settings while maintaining performance. More importantly, in contrast to previously proposed trainable calibration objectives, ESD does not contain any internal hyperparameters, which significantly reduces the total computational cost for training. This reduction in cost is more prominent as the complexity of the model and dataset increases, making ESD a more viable calibration objective option.

## ACKNOWLEDGEMENT

This work was supported by Institute of Information & communications Technology Planning & Evaluation (IITP) grant funded by the Korea government(MSIT) (No.2022-0-00184, Development and Study of AI Technologies to Inexpensively Conform to Evolving Policy on Ethics), and Institute for Information & communications Technology Promotion(IITP) grant funded by the Korea government(MSIT) (No. 2021-0-01381, Development of Causal AI through Video Understanding and Reinforcement Learning, and Its Applications to Real Environments).

## REFERENCES

- Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural network. *International conference on machine learning*, pp. 1613–1622. PMLR, 2015.
- Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Prasoon Goyal, Lawrence D. Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, Xin Zhang, Jake Zhao, and Karol Zieba. End to end learning for self-driving cars. *CoRR abs/1604.07316*, 2016. URL <http://arxiv.org/abs/1604.07316>
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. A large annotated corpus for learning natural language inference. *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Association for Computational Linguistics, 2015.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. *2009 IEEE Conference on Computer Vision and Pattern Recognition* pp. 248–255, 2009. doi: 10.1109/CVPR.2009.5206848.
- Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE Signal Processing Magazine* 29(6):141–142, 2012.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)* 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL: <https://aclanthology.org/N19-1423>
- Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. *International conference on machine learning*, pp. 1050–1059. PMLR, 2016.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural networks. In Doina Precup and Yee Whye Teh (eds), *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 1321–1330. PMLR, 06–11 Aug 2017. URL <https://proceedings.mlr.press/v70/guo17a.html>
- Kartik Gupta, Amir Rahimi, Thalaiyasingam Ajanthan, Thomas Mensink, Cristian Sminchisescu, and Richard Hartley. Calibration of neural networks using splines. *International Conference on Learning Representations* 2021. URL <https://openreview.net/forum?id=eQe8DEWNN2W>
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* 770–778, 2016. doi: 10.1109/CVPR.2016.90.
- Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *International Conference on Learning Representations* 2019. URL <https://openreview.net/forum?id=HJz6tiCqYm>

- Archit Karandikar, Nicholas Cain, Dustin Tran, Balaji Lakshminarayanan, Jonathon Shlens, Michael Curtis Mozer, and Rebecca Roelofs. Soft calibration objectives for neural networks. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds), *Advances in Neural Information Processing Systems*, 2021. URL <https://openreview.net/forum?id=-tVD13hOsQ3>.
- Junyeong Kim and Chang D. Yoo. Deep partial person re-identification via attention model. In 2017 IEEE International Conference on Image Processing (ICIP), pp. 3425–3429, 2017. doi: 10.1109/ICIP.2017.8296918.
- ByungSoo Ko, Ho-Jin Choi, Chansol Hong, Jong-Hwan Kim, Oh Chul Kwon, and Chang D. Yoo. Neural network-based autonomous navigation for a homecare mobile robot. In 2017 IEEE International Conference on Big Data and Smart Computing (BigComp), pp. 403–406, 2017. doi: 10.1109/BIGCOMP.2017.7881744.
- Simon Kocbek, Primoz Kocbek, Leona Cilar, and Gregor Stiglic. Local interpretability of calibrated prediction models: A case of type 2 diabetes mellitus screening test, 2020. <https://arxiv.org/abs/2006.13815>.
- Ranganath Krishnan and Omesh Tickoo. Improving model calibration with accuracy versus uncertainty optimization. In *Proceedings of the 34th International Conference on Neural Information Processing Systems, NIPS'20*, Red Hook, NY, USA, 2020. Curran Associates Inc. ISBN 9781713829546.
- Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-10 (canadian institute for advanced research). a. URL <http://www.cs.toronto.edu/~kriz/cifar.html>.
- Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-100 (canadian institute for advanced research). b. URL <http://www.cs.toronto.edu/~kriz/cifar.html>.
- Aviral Kumar, Sunita Sarawagi, and Ujjwal Jain. Trainable calibration measures for neural networks from kernel mean embeddings. In Jennifer Dy and Andreas Krause (eds), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 2805–2814. PMLR, 10–15 Jul 2018. URL <https://proceedings.mlr.press/v80/kumar18a.html>.
- Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. *Advances in neural information processing systems* 30, 2017.
- Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. doi: 10.1109/5.726791.
- Jing Lei, Max G'Sell, Alessandro Rinaldo, Ryan J Tibshirani, and Larry Wasserman. Distribution-free predictive inference for regression. *Journal of the American Statistical Association*, 113(523): 1094–1111, 2018.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=Bkg6RiCqY7>.
- Bill MacCartney and Christopher D. Manning. Modeling semantic containment and exclusion in natural language inference. *Proceedings of the 22nd International Conference on Computational Linguistics (Coling 2008)*, pp. 521–528, Manchester, UK, August 2008. Coling 2008 Organizing Committee. URL <https://aclanthology.org/C08-1066>.
- Jishnu Mukhoti, Viveka Kulharia, Amartya Sanyal, Stuart Golodetz, Philip H. S. Torr, and Puneet K. Dokania. Calibrating deep neural networks using focal loss. *Proceedings of the 34th International Conference on Neural Information Processing Systems, NIPS'20*, Red Hook, NY, USA, 2020. Curran Associates Inc. ISBN 9781713829546.

- Mahdi Pakdaman Naeni, Gregory F Cooper, and Milos Hauskrecht. Obtaining well calibrated probabilities using bayesian binning. *AAAI*, pp. 2901–2907, 2015.
- Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. Adversarial NLI: A new benchmark for natural language understanding. *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 4885–4901, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.441. URL: <https://aclanthology.org/2020.acl-main.441>
- Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, David Sculley, Sebastian Nowozin, Joshua Dillon, Balaji Lakshminarayanan, and Jasper Snoek. Can you trust your model's uncertainty? evaluating predictive uncertainty under dataset shift. *Advances in neural information processing systems*, 32, 2019.
- Mahdi Pakdaman Naeni, Gregory Cooper, and Milos Hauskrecht. Obtaining well calibrated probabilities using bayesian binning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 28(1), Feb. 2015. URL: <https://ojs.aaai.org/index.php/AAAI/article/view/9602>
- Trung Xuan Pham, Jin Woong Choi, Rusty John Lloyd Mina, Thanh Xuan Nguyen, Sultan Rizky Madjid, and Chang D. Yoo. Lad: A hybrid deep learning system for benign paroxysmal positional vertigo disorders diagnosis. *IEEE Access*, 10:113995–114007, 2022. doi: 10.1109/ACCESS.2022.3215625.
- John C. Platt. Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. In *ADVANCES IN LARGE MARGIN CLASSIFIERS*, pp. 61–74. MIT Press, 1999.
- A. W. van der Vaart. *Asymptotic Statistics Cambridge Series in Statistical and Probabilistic Mathematics*. Cambridge University Press, 1998. doi: 10.1017/CBO9780511802256.
- Vladimir Vovk, Alexander Gammerman, and Glenn Shafer. *Algorithmic learning in a random world*. Springer Science & Business Media, 2005.
- Sunjae Yoon, Dahyun Kim, Ji Woo Hong, Junyeong Kim, and Chang D. Yoo. Dual-scale doppler attention for human identification. *Sensors*, 22(17), 2022. ISSN 1424-8220. doi: 10.3390/s22176363. URL: <https://www.mdpi.com/1424-8220/22/17/6363>
- Bianca Zadrozny and Charles Elkan. Obtaining calibrated probability estimates from decision trees and naive bayesian classifiers. *Proceedings of the Eighteenth International Conference on Machine Learning (ICML '01)*, pp. 609–616, San Francisco, CA, USA, 2001. Morgan Kaufmann Publishers Inc. ISBN 1558607781.
- Bianca Zadrozny and Charles Elkan. Transforming classifier scores into accurate multiclass probability estimates. In *Proceedings of the Eighth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '02)*, pp. 694–699, New York, NY, USA, 2002. Association for Computing Machinery. ISBN 158113567X. doi: 10.1145/775047.775151. URL: <https://doi.org/10.1145/775047.775151>

## A VISUAL INTUITION OF EXPECTED SQUARED DIFFERENCE (ESD)

From Eq. (7),

$$\begin{aligned} d_k(\cdot) &= \int_0^Z P(Y_k = 1; Z_k = z_k) - z_k P(Z_k = z_k) dz_k \\ &= \int E_{Z_k; Y_k} [I(Z_k \leq z_k) - (Y_k = 1) I(Z_k \leq z_k)] dz_k \end{aligned}$$

This can be viewed as the difference between two quantities:

$$\begin{aligned} \text{Cumulative Accuracy} &= \int_0^Z P(Y_k = 1; Z_k = z_k) dz_k \\ &= E_{Z_k; Y_k} [I(Z_k \leq Z; Y_k = 1)] \\ &= \frac{1}{N} \sum_{i=1}^N I(z_{k;i} \leq Z; Y_{k;i} = 1) \end{aligned}$$

$$\begin{aligned} \text{Cumulative Confidence} &= \int_0^Z z_k P(Z_k = z_k) dz_k \\ &= E_{Z_k; Y_k} [Z_k I(Z_k \leq Z)] \\ &= \frac{1}{N} \sum_{i=1}^N z_{k;i} I(z_{k;i} \leq Z) \end{aligned}$$

Due to the close relationship between  $d_k(\cdot)$  and the calibration of a neural network, the average squared difference  $E_{Z_k} [d_k(Z_k^0)^2]$  between the cumulative accuracy and confidence closely corresponds with the calibration of a network (Figure 4). That is, the average squared difference between the cumulative accuracy and confidence is larger for an uncalibrated network compared to a calibrated network. Jointly training with our proposed Expected Squared Difference (ESD) as an auxiliary loss tries to minimize this squared difference between the two curves during training on average, thus achieving a better calibrated model.

Figure 4: Visual intuition plot showing the cumulative confidence and cumulative accuracy with varying quantile scores of prediction confidence for an uncalibrated (left) and calibrated (right) network. The uncalibrated network was obtained by training Resnet34 on CIFAR100 with NLL, and the calibrated network was acquired by temperature scaling on the aforementioned trained network.

## B PROOF THAT ESD IS AN UNBIASED AND CONSISTENT ESTIMATOR

In Theorem 2, we prove that ESD is an unbiased estimator, that is taking the expectation of ESD is equal to the true Expected Squared Difference.

Theorem 2 ESD is an unbiased estimator, i.e.  $E_{Z_k; Y_k} [\text{ESD}] = E_{Z_k^0} [d_k^2(Z_k^0)]$  where  $d_k(Z_k^0) = \int E_{Z_k; Y_k} [I(Z_k \leq Z_k^0) - (Y_k = 1) I(Z_k \leq Z_k^0)] dz_k$ .

Proof. Let  $Z_k = (Z_{k;1}, \dots, Z_{k;n})$ ,  $Y_k = (Y_{k;1}, \dots, Y_{k;n})$  and  $G_i = g^2 \frac{S_{g_i}^2}{N-1}$ . We have that

$$\begin{aligned} E_{Z_k; Y_k}[\text{ESD}] &= \frac{1}{N} \sum_{i=1}^N E_{Z_k; Y_k}[G_i] \quad (\text{by linearity of expectation}) \\ &= \frac{1}{N} \sum_{i=1}^N E_{Z_{k;i}; Y_{k;i}}[g_i^2] \quad (\text{by lemma 2.1}) \\ &= \frac{1}{N} \sum_{i=1}^N E_{Z_k^0}[d_k^2(Z_k^0)] \quad (\text{since } E_{Z_{k;i}; Y_{k;i}}[g_i^2] = E_{Z_k^0}[d_k^2(Z_k^0)]) \\ &= E_{Z_k^0}[d_k^2(Z_k^0)]: \end{aligned}$$

□

Lemma 2.1  $E_{Z_{k;i}; Y_{k;i}}[G_i] = \frac{S_{g_i}^2}{N-1}$  where  $S_{g_i}^2 = E_{Z_{k;i}; Y_{k;i}}[(Z_{k;i} - Z_{k;i})^2 | (Y_{k;i} = 1) - z_{k;i}]$ :

Proof. Since samples are i.i.d., for a fixed  $i$  it holds that:

$$\begin{aligned} E_{Z_{k;i}; Y_{k;i}}[g_i] &= g_i \\ E_{Z_{k;i}; Y_{k;i}}\left[\frac{S_{g_i}^2}{N-1}\right] &= \frac{S_{g_i}^2}{N-1} \quad (\text{where } \text{Var}[g_i] = \frac{S_{g_i}^2}{N-1}) \\ \text{Var}[g_i] &= E_{Z_{k;i}; Y_{k;i}}[g_i^2] - \frac{S_{g_i}^2}{N-1}: \end{aligned} \tag{12}$$

Shifting variables around, we get

$$\begin{aligned} \frac{S_{g_i}^2}{N-1} &= E_{Z_{k;i}; Y_{k;i}}[g_i^2] - \text{Var}[g_i] \\ &= E_{Z_{k;i}; Y_{k;i}}[g_i^2] - \frac{S_{g_i}^2}{N-1} \quad (\text{since } \text{Var}[g_i] = \frac{S_{g_i}^2}{N-1}) \\ &= E_{Z_{k;i}; Y_{k;i}}[G_i]: \end{aligned}$$

□

In Theorem 3, we prove that ESD is a consistent estimator, which means that ESD converges in probability to the true Expected Squared Difference.

Theorem 3 ESD is a consistent estimator, i.e.  $E_{Z_k^0}[\text{ESD}] \xrightarrow{P} E_{Z_k^0}[d_k^2(Z_k^0)]:$

Proof. Since ESD is an unbiased estimator, it is sufficient to prove  $\lim_{n \rightarrow \infty} \text{Var}[\text{ESD}] = 0$ . For simplicity, we use where  $\sum_{i,j=1}^N = \sum_{i=1}^N \sum_{j=1}^N$  and  $E_{Z_k; Y_k}[G_i] = E_{Z_k^0}[d_k^2(Z_k^0)] \delta_i$ ;

$$\text{Var}[\text{ESD}] = \frac{1}{N^2} \sum_{i=1}^N \text{Var}[G_i] + \frac{1}{N^2} \sum_{\substack{i,j=1 \\ i \neq j}}^N \text{Cov}[G_i; G_j];$$

$$\begin{aligned} \text{Since } E_{Z_k; Y_k}[G_i] &= E_{Z_k^0}[d_k^2(Z_k^0)] \delta_i; \\ &= \frac{1}{N^2} \sum_{i=1}^N \text{Var}[G_i] + \frac{1}{N^2} \sum_{\substack{i,j=1 \\ i \neq j}}^N (E_{Z_k; Y_k}[G_i G_j] - E_{Z_k^0}^2[d_k^2(Z_k^0)]) \\ &= \frac{1}{N^2} \sum_{i=1}^N \text{Var}[G_i] + \frac{1}{N^2} \sum_{\substack{i,j=1 \\ i \neq j}}^N (E_{Z_k; Y_k}[g_1^2 g_2^2] - E_{Z_k^0}^2[d_k^2(Z_k^0)]) \quad (\text{by lemma 3.1 and 3.2}) \\ &\quad + \frac{1}{N^2} \sum_{\substack{i,j=1 \\ i \neq j}}^N \frac{4}{(N-2)^2} \\ &= \frac{2}{N} + \frac{4N(N-1)}{N^2(N-2)^2} + \frac{N(N-1)}{N^2} (E_{Z_k; Y_k}[g_1^2 g_2^2] - E_{Z_k^0}^2[d_k^2(Z_k^0)]); \quad (\text{by lemma 3.3}) \end{aligned}$$

Thus, by using lemma 3.4 and the non-negativeness of variance,

$$\lim_{n \rightarrow \infty} \text{Var}[\text{ESD}] = \lim_{n \rightarrow \infty} E_{Z_k; Y_k}[g_1^2 g_2^2] - E_{Z_k^0}^2[d_k^2(Z_k^0)] = 0;$$

Consequently  $\lim_{n \rightarrow \infty} \text{Var}[\text{ESD}] = 0$ : □

Lemma 3.1  $\text{Var}[G_i] \leq 2 < 1$ .

Proof.

$$\begin{aligned} \text{Var}[G_i] &= E_{Z_k; Y_k}[G_i^2] - E_{Z_k; Y_k}^2[G_i] \\ &= E_{Z_k; Y_k}[G_i^2] - \frac{S_{g_i}^2}{N-1} \\ &= E_{Z_k; Y_k}[4g_1^2 g_2^2] - \frac{S_{g_i}^2}{N-1} \\ &= E_{Z_k; Y_k}[4g_1^2 g_2^2] + \frac{S_{g_i}^2}{N-1} \cdot \frac{2}{N-2} \end{aligned}$$

By lemma 3.2  $E_{Z_k; Y_k}[g_1^2 g_2^2] \leq 1$  and  $E_{Z_k; Y_k} \frac{S_{g_i}^2}{N-1} \leq \frac{2}{N-2}$ . Thus,

$$E_{Z_k; Y_k}[g_1^2 g_2^2] + \frac{S_{g_i}^2}{N-1} \leq 1 + \frac{2}{N-2} < 2;$$

□

Lemma 3.2  $E_{Z_k; Y_k} \frac{S_{g_i}^2}{N-1} \leq \frac{2}{N-2}$  and  $E_{Z_k; Y_k}[g_1^2 g_2^2] \leq 1$ :

Proof. By the triangle inequality,  $\|g_i - g_j\| = \left\| \frac{1}{N-1} \sum_{m \in i} g_{im} - \frac{1}{N-1} \sum_{m \in j} g_{jm} \right\| \leq 1$ , since  $\|g_{im} - g_{jm}\| \leq 1$ .

Furthermore, we have

$$\begin{aligned} \frac{S_{g_i}^2}{N-1} &= \frac{1}{(N-1)^2} \sum_{m \in i} g_{im}^2 + \frac{1}{(N-1)^2(N-2)} \left( \sum_{m \in i} g_{im} \right)^2 \\ &\quad - \frac{1}{(N-1)^2} \sum_{m \in i} \|g_{im}\|^2 + \frac{1}{(N-1)^2(N-2)} \left( \sum_{m \in i} \|g_{im}\| \right)^2 \\ &\quad - \frac{1}{(N-1)^2} \sum_{m \in i} 1 + \frac{1}{(N-1)^2(N-2)} \left( \sum_{m \in i} 1 \right)^2. \end{aligned}$$

The last inequality is by  $\|g_{im}\| \leq 1$ , and the direct calculation implies that

$$\frac{S_{g_i}^2}{N-1} \leq \frac{1}{N-1} + \frac{1}{N-2} - \frac{2}{N-2}.$$

□

Lemma 3.3  $E_{Z_k; Y_k} [g_1^2 g_2^2] = E_{Z_k; Y_k} [g_1^2 g_2^2] \delta_{ij}$  where  $i \in j$ :

Proof.

$$\begin{aligned} \sum_{\substack{m, n=1 \\ f, m, n \in \mathcal{F} \cap i, j, g}}^N g_{im} g_{jn} &= \sum_{\substack{m, n=1 \\ f, m, n \in \mathcal{F} \cap i, j, g, m \in n}}^N g_{im} g_{jn} + \sum_{\substack{m, n=1 \\ n \in \mathcal{F} \cap i, j, g}}^N g_{ij} g_{jn} + \sum_{\substack{m, n=1 \\ m \in \mathcal{F} \cap i, j, g}}^N g_{im} g_{ji} \\ &\quad + \sum_{\substack{m, n=1 \\ m \in \mathcal{F} \cap i, j, g}}^N g_{im} g_{im} + g_{ij} g_{ji} \quad (\text{since these terms are disjoint terms}) \\ &\stackrel{d}{=} \sum_{\substack{m^0, n^0=1 \\ f, m^0, n^0 \in \mathcal{F} \cap 1; 2g, m^0 \in n^0}}^N g_{1m^0} g_{2n^0} + \sum_{\substack{m^0, n^0=1 \\ n^0 \in \mathcal{F} \cap 1; 2g}}^N g_{12} g_{2n^0} + \sum_{\substack{m^0, n^0=1 \\ m^0 \in \mathcal{F} \cap 1; 2g}}^N g_{1m^0} g_{21} \\ &\quad + \sum_{\substack{m^0, n^0=1 \\ m^0 \in \mathcal{F} \cap 1; 2g}}^N g_{1m^0} g_{2m^0} + g_{12} g_{21} \\ &= \sum_{\substack{m^0, n^0=1 \\ m^0 \in \mathcal{F} \cap 1; n^0 \in \mathcal{F} \cap 2}}^N g_{1m^0} g_{2n^0}. \end{aligned}$$

Proof of claim above:

As the summation can be divided into  $v$  distinct cases,



Let  $i$  and  $j$  be arbitrary such that  $i \neq j$ ,

Case 1  $f = m; n \in f; j \in g$  and  $m \in n$ ,

$$\begin{aligned} g_m g_n &= I(Z_{k:m} = Z_{k:i}) [I(Y_{k:m} = 1) - Z_{k:m}] I(Z_{k:n} = Z_{k:j}) [I(Y_{k:n} = 1) - Z_{k:n}] \\ &= h_1(Z_{k:i}; Z_{k:m}; Z_{k:j}; Z_{k:n}; Y_{k:m}; Y_{k:n}) \\ &\stackrel{d}{=} h_1(Z_{k;1}; Z_{k;m^0}; Z_{k;2}; Z_{k;n^0}; Y_{k;m^0}; Y_{k;n^0}): \quad \text{where } m^0, n^0 \in f; j; 2g \text{ and } m^0 \in n^0 \end{aligned}$$

Case 2  $m = j$  and  $n \in f; j \in g$ ,

$$\begin{aligned} g_j g_n &= I(Z_{k:j} = Z_{k:j}) [I(Y_{k:j} = 1) - Z_{k:j}] I(Z_{k:n} = Z_{k:j}) [I(Y_{k:n} = 1) - Z_{k:n}] \\ &= h_2(Z_{k:i}; Z_{k:j}; Z_{k:n}; Y_{k;j}; Y_{k:n}) \\ &\stackrel{d}{=} h_2(Z_{k;1}; Z_{k;2}; Z_{k;n^0}; Y_{k;2}; Y_{k;n^0}): \quad \text{where } m^0 = 2 \text{ and } n^0 \in f; j; 2g \end{aligned}$$

Case 3  $n = i$  and  $m \in f; j \in g$ ,

$$\begin{aligned} g_m g_i &= I(Z_{k:m} = Z_{k:i}) [I(Y_{k:m} = 1) - Z_{k:m}] I(Z_{k:i} = Z_{k:j}) [I(Y_{k:i} = 1) - Z_{k:i}] \\ &= h_3(Z_{k:i}; Z_{k:m}; Z_{k:j}; Y_{k;m}; Y_{k;i}) \\ &\stackrel{d}{=} h_3(Z_{k;1}; Z_{k;m^0}; Z_{k;2}; Y_{k;m^0}; Y_{k;1}): \quad \text{where } n^0 = 1 \text{ and } m^0 \in f; j; 2g \end{aligned}$$

Case 4  $m = n \in f; j \in g$ ,

$$\begin{aligned} g_m g_m &= I(Z_{k:m} = Z_{k:i}) [I(Y_{k:m} = 1) - Z_{k:m}] I(Z_{k:m} = Z_{k:j}) [I(Y_{k:m} = 1) - Z_{k:m}] \\ &= h_4(Z_{k:i}; Z_{k:m}; Z_{k:j}; Y_{k;m}) \\ &\stackrel{d}{=} h_4(Z_{k;1}; Z_{k;m^0}; Z_{k;2}; Y_{k;m^0}): \quad \text{where } m^0 = n^0 \in f; j; 2g \end{aligned}$$

Case 5  $m = j$  and  $n = i$ ,

$$\begin{aligned} g_j g_i &= I(Z_{k:j} = Z_{k:i}) [I(Y_{k:j} = 1) - Z_{k:j}] I(Z_{k:i} = Z_{k:j}) [I(Y_{k:i} = 1) - Z_{k:i}] \\ &= h_5(Z_{k:i}; Z_{k:j}; Y_{k;i}; Y_{k;j}) \\ &\stackrel{d}{=} h_5(Z_{k;1}; Z_{k;2}; Y_{k;1}; Y_{k;2}): \quad \text{where } m^0 = 2 \text{ and } n^0 = 1 \end{aligned}$$

This follows from the fact that  $Z_{k,i}$ 's and  $Y_{k,i}$ 's are i.i.d. random variables. Therefore,

$$\begin{aligned} \left( \prod_{\substack{m:n=1 \\ m \in i; n \in j}}^N g_m g_n \right)^2 &\stackrel{d}{=} \left( \prod_{\substack{m^0:n^0=1 \\ m^0 \in 1; n^0 \in 2}}^N g_{1m^0} g_{2n^0} \right)^2 \\ g_i^2 g_j^2 &\stackrel{d}{=} g_1^2 g_2^2 \\ E_{Z_k; Y_k} [g_i^2 g_j^2] &= E_{Z_k; Y_k} [g_1^2 g_2^2]: \end{aligned}$$

□

Lemma 3.4  $\lim_{N \rightarrow \infty} E_{Z_k; Y_k} [g_1^2 g_2^2] = E_{Z_k^0} [d_k^2(Z_k^0)]$ .

Proof. For arbitrary  $\epsilon > 0$ ,  $g_i \xrightarrow{P} 1$  by the strong law of large number. Thus,  $g_1 \xrightarrow{P} 1$  and  $g_2 \xrightarrow{P} 2$ . Therefore, since marginal convergence in probability implies joint convergence in probability (Vaart, 1998),

$$(g_1; g_2) \xrightarrow{P} (1; 2):$$

By continuous mapping theorem, since  $f(x; y) = x^2 y^2$  is a continuous function,

$$g_1^2 g_2^2 \xrightarrow{P} 1^2 2^2:$$

Additionally, since  $g_1^2 g_2^2 \leq 1 \cdot 8N$ , it is uniformly bounded and thus uniformly integrable. Combining this with the fact that it converges in probability,

$$\lim_{N \rightarrow \infty} E_{Z_k; Y_k} [g_1^2 g_2^2] = E_{Z_k; Y_k} [1^2 2^2]:$$

Thus,

$$\lim_{N \uparrow} E_{Z_k; Y_k} [g_1^2 g_2^2] = E_{Z_k; Y_k} [ \begin{matrix} 2 & 2 \\ 1 & 2 \end{matrix} ] = E_{Z_k; Y_k} [ \begin{matrix} 2 \\ 1 \end{matrix} ] E_{Z_k; Y_k} [ \begin{matrix} 2 \\ 2 \end{matrix} ] = E_{Z_k^0}^2 [d_k^2(Z_k^0)]:$$

□

## C EXPECTED SQUARED DIFFERENCE (ESD) PSEUDOCODE

In this section, we provide a pseudocode for calculating the Expected Squared Difference (ESD) for a given batch output.

---

### Algorithm 1: Pytorch-like Pseudocode: Expected Squared Difference (ESD)

---

```
# n: number of samples in a mini-batch.
# confidence: 1D tensor of n elements containing max softmax
# outputs of each sample from a neural network.
# correct: 1D tensor of n elements containing 1/0 corresponding
# to the correctness of each sample.
def ESD_loss(n, confidence, correct):
    # compute the difference between confidence and correctness.
    diff = correct.float() - confidence

    # Prepare the split between inner and outer expectation
    # estimation.
    split = torch.ones(n, n) - torch.eye(n)

    # compute the inner expectation estimation.
    confidence_mat = confidence.expand(n, n)
    ineq = torch.le(confidence_mat, confidence_mat.T).float()
    diff_mat = diff.view(1, n).expand(n, n)
    x_mat = torch.mul(diff_mat, ineq) * split
    mean_row = torch.sum(x_mat, dim = 1) / (n - 1)
    x_mat_squared = torch.mul(x_mat, x_mat)
    var = 1 / (n - 2) * torch.sum(x_mat_squared, dim = 1) - (n - 1) / (n - 2)
    torch.mul(mean_row, mean_row)

    # compute the outer expectation estimation.
    d_k_sq_vector = torch.mul(mean_row, mean_row) - var / (n - 1)
    ESD = torch.sum(d_k_sq_vector) / n
    return ESD
```

---

## D UNBIASED ESTIMATORS AND ITS GRADIENT

Since ESD is an unbiased estimator,

$$\begin{aligned} E_{Z_k; Y_k} [\text{ESD}] &= E_{Z_k^0} [d_k^2(Z_k^0)] \\ r E_{Z_k^0} [d_k^2(Z_k^0)] &= r E_{Z_k; Y_k} [\text{ESD}] \\ &= r E_{X; Y} [\text{ESD}] \quad (\text{by law of unconscious statistician}) \\ &= E_{X; Y} [r \text{ESD}]: \quad (\text{since } (X; Y) \text{ are independent of } \end{aligned}$$

Thus, the gradient of the unbiased estimator is on average the gradient of the desired metric.

## E ADDITIONAL INFORMATION ON THE CHOICE OF LAMBDA RANGE

To validate the choice of grid range in our experimental settings, we plot the variations in accuracy with respect to increasing for all methods (Figure 5). We see that our choice contains points

Figure 5: Accuracy plot with respect to varying values  $\alpha$  across different datasets and models trained with MMCE, SB-ECE, and ESD. The threshold accuracy represents the value 1.5% below the baseline accuracy, which was used as the model selection criterion as stated in section 5.2.

within the acceptable range of accuracy (i.e., within 1.5% degradation in accuracy compared to baseline). However, after a particular value, which is different on a per method per dataset basis, the accuracy decreases consistently with increasing  $\alpha$ . As such, the grid chosen is suitable for the experiments conducted.

## F TRAINABLE CALIBRATION MEASURES UNDER DISTRIBUTION SHIFT

Figure 6: Accuracy and ECE plot for varying intensities of distribution shifts in (a) CIFAR 10-C and (b) CIFAR 100-C across models trained using NLL as well as those jointly trained with an auxiliary calibration objective (i.e., NLL+MMCE, NLL+SB-ECE, NLL+ESD) and their performance after post-processing with temperature scaling (i.e., NLL+T, NLL+MMCE+T, NLL+SB-ECE+T, NLL+ESD+T).

In Figure 6, we evaluate the performance of calibration methods under distribution shift benchmark datasets, CIFAR 10-C and CIFAR 100-C, introduced in Hendrycks & Dietterich (2019). For models jointly trained with an auxiliary calibration objective (NLL+MMCE, NLL+SB-ECE, NLL+ESD), we observe that they are more robust to distribution shifts when compared to those solely trained with NLL. In addition, prior work (Ovadia et al., 2019), has shown that the calibration performance of a neural network after temperature scaling may be significantly reduced under distribution shifts. Our results on CIFAR 10-C and CIFAR 100-C imply a similar trend. Stacking calibration during training methods with temperature scaling (NLL+MMCE+T, NLL+SB-ECE+T, NLL+ESD+T), we observe that they perform comparable to temperature scaled models trained with NLL (NLL+T) with the exception of ESD, where it performs marginally better. As such, training with ESD could potentially improve a model's robustness after temperature scaling regarding distribution shifts and mitigate this issue.