

PEACOCK: MULTI-OBJECTIVE OPTIMIZATION FOR DEEP NEURAL NETWORK CALIBRATION

Anonymous authors

Paper under double-blind review

ABSTRACT

The rapid adoption of deep neural networks underscores an urgent need for models to be safe, trustworthy and well-calibrated. Despite recent advancements in network calibration, the optimal combination of techniques remains relatively unexplored. By framing the task as a multi-objective optimization problem, we demonstrate that combining state-of-the-art methods can further boost calibration performance. We feature a total of seven state-of-the-art calibration algorithms and provide both theoretical and empirical motivation for their equal and weighted importance unification. We conduct experiments on both in-distribution and out-of-distribution computer vision and natural language benchmarks, investigating the speeds and contributions of different components. Our code is available anonymously at: <https://anonymous.4open.science/r/Peacock-1CE8>.

1 INTRODUCTION

Key requirements for the safe deployment of neural networks include multiple desirable qualities, such as high accuracy, fast training speeds and trustworthy predictions. While the recent successes in deep learning have increased the use of complex neural networks, a common observation is that deep models tend to be miscalibrated, exhibiting either under- or over-confident predictions (Guo et al., 2017).

Miscalibration can be particularly dangerous for high-stakes, safety-critical tasks such as medical prognosis (Esteva et al., 2017; Bandi et al., 2019), object-detection (Munir et al., 2023a;b; Liu et al., 2024), AI fairness and decision-making (Pleiss et al., 2017; Corvelo Benz & Rodriguez, 2023). Such tasks demand reliable decision-making algorithms, necessitating accurate confidence estimates that reflect a model’s uncertainty (Jiang et al., 2018; Kendall & Gal, 2017). Specifically, calibration ensures that a model’s predicted confidences align with its actual correctness. For instance, if a model assigns 0.9 confidence to a set of 100 samples, we should expect the model to be correct for 90 instances only.

Modern neural networks must not only remain well-calibrated in-distribution (ID), but also display invariance properties and remain robustly calibrated against out-of-distribution (OOD) shifts (Wald et al., 2021). This is crucial for real-world deployment, where models must generalize well and express uncertainty when handling unseen inputs (see Fig. 2b). For instance, OOD shifts in computer vision might involve changes in saturation and illumination, while in natural language tasks, they can arise from differences in syntax or spelling mistakes (Zhang et al., 2023).

While most calibration techniques tend to outperform the vanilla cross-entropy (CE) loss, they each tackle radically different issues, enabling independent performance boosts through different approaches. Each of these techniques exhibit varying trade-offs in ID/OOD performance (see Fig. 1 showing the 100% - ECE%)

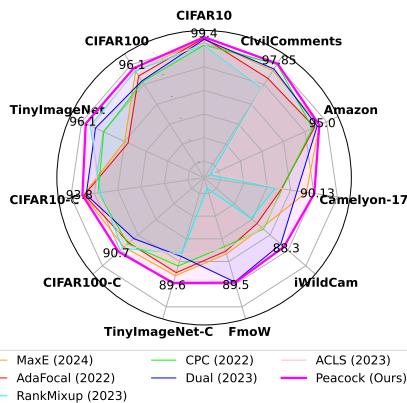
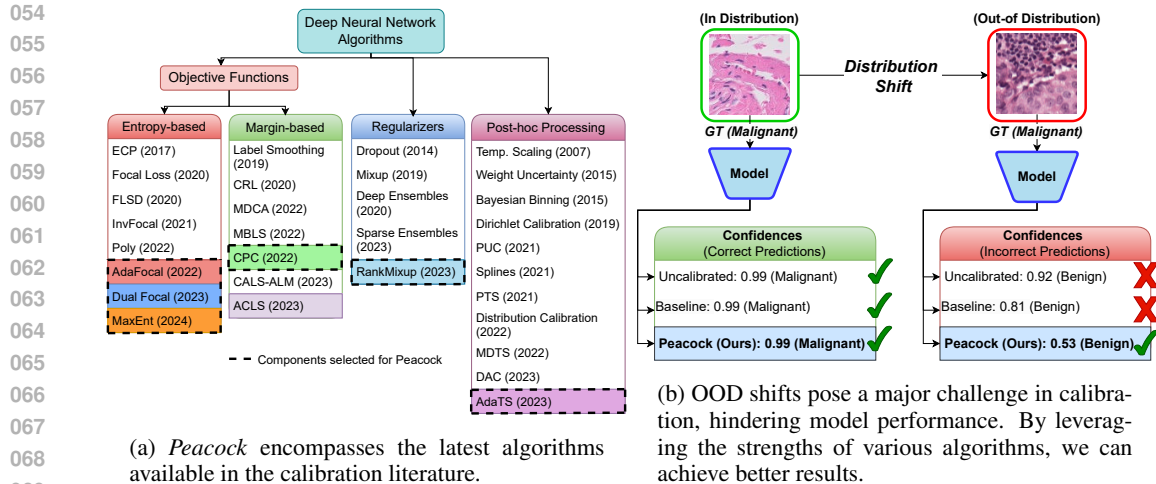


Figure 1: By unifying multiple calibration algorithms, *Peacock* outperforms individual methods, achieving state-of-the-art ID and OOD calibration. (Bigger is better)



070 Figure 2: Motivated by recent contributions, we propose *Peacock* for calibration under OOD shifts.

073 and computational complexity, making it difficult to pinpoint a clear winner (Cheng & Vasconcelos, 2022). Although it is theoretically possible for many calibration methods to be combined, the question remains on *which* of and *how* can these methods be fruitfully integrated together?

076 In this paper, our main goal is to *unify* different proposed calibration approaches into a single calibration framework named *Peacock*. Our claim is that by jointly optimizing multiple calibration objectives, performance boosts can be achieved for ID + OOD classification tasks. Theoretically, we demonstrate that *Peacock*'s calibration errors will always be bounded by the average calibration errors of all its components. We also propose a novel weighted importance form of *Peacock* that is fast and effective in balancing the contributions of different components. Apart from presenting *Peacock*, this paper is also doubly positioned as a literature survey of recently published calibration strategies (see Fig. 2a).

084 To demonstrate the combined efficacy of different calibration methods, we revisit a total of seven different SOTA calibration baselines and pick the final six to be integrated into *Peacock*. We further evaluate both ID and OOD performances of *Peacock* on popular synthetic and in-the-wild computer vision and natural language tasks. Our contributions can be summarized into the following points.

- 088 • **Peacock:** We present *Peacock*, a fully integrated multi-objective framework for deep neural network calibration.
- 089 • **Equal and Weighted Importance:** We motivate *Peacock* with theoretical guarantees and propose novel ways to weight the contributions of different calibration components.
- 090 • **Review of Literature:** This paper additionally serves as a condensed survey¹ of all existing algorithms proposed in the calibration literature.
- 091 • **Evaluation and Analysis:** We evaluate across popular synthetic and in-the-wild OOD vision and text benchmarks, empirically analyzing the speeds, effects and contributions of different components.

098 2 RELATED WORK

100 **Multi-Objective Optimization** The primary goal of multi-objective optimization (MOO) is to simultaneously optimize multiple loss² terms $\{\mathcal{L}_1, \dots, \mathcal{L}_A\}$ on a single model. As these differing loss functions may possibly contrast and conflict with each other (Sener & Koltun, 2018). Obtaining a way to balance/weight these losses during optimization is of great interest, since learned representations between losses can be shared (Caruana, 1997; Zamir et al., 2018), with the added benefit

106 ¹ To the best of our ability, we cover all published works and further discuss them in Appendix A

107 ² For clarity, the term **loss** loosely refers to auxiliary functions, whilst **objective function** represents the final learning goal during model training.

of avoiding model redundancy in the form of large ensembles (Dosovitskiy & Djolonga, 2020). Prevalent methods include grid-search tuning, which can be cumbersome and fixing predefined weights during training may not guarantee optimal performance (Groenendijk et al., 2021). Other approaches include learning the hardest task first (Guo et al., 2018), self-paced learning (Li et al., 2017), aleatoric uncertainty estimation (Kendall & Gal, 2017; Kendall et al., 2018), gradient normalization (Chen et al., 2018), pareto frontiers (Sener & Koltun, 2018; Lin et al., 2019; Xiao et al., 2023; Liu et al., 2021) and co-efficient of variations (Groenendijk et al., 2021). Additionally, we refer readers to (Zhang & Yang, 2017; Gong et al., 2019) for a comprehensive review and comparisons on multi-objective methods. In this work, our focus is directed towards gradient-based multi-objective optimization for balancing calibration loss terms.

Deep Uncertainty Calibration In Fig. 2a, we provide an overview of recent calibration algorithms and metrics. Examples of these algorithms include (1) *Entropy-based* methods that control the entropy of the model (Mukhoti et al., 2020; Wang et al., 2021; Leng et al., 2022; Ghosh et al., 2022; Tao et al., 2023; Neo et al., 2024). (2) *Margin-based* methods that directly limit model confidences (Hebbalaguppe et al., 2022; Liu et al., 2022a; Cheng & Vasconcelos, 2022; Liu et al., 2023a;b) (3) *Regularizers* that augment the training inputs or model (Zhang et al., 2020; Sapkota et al., 2023; Noh et al., 2023). (4) *Post-hoc processing*, which requires tuning the model on a hold-out validation set in order to scale predictions (Wenger et al., 2020; Tomani et al., 2021; Gupta et al., 2021; Tomani et al., 2022; Kuleshov & Deshpande, 2022; Gruber & Buettner, 2022; Yu et al., 2022; Tomani et al., 2023; Joy et al., 2023). Calibration metrics include binning and binning-free approaches (Gupta et al., 2021; Roelofs et al., 2022; Yang et al., 2023; Xiong et al., 2023). Additionally, we refer readers to Appendix A.5 for a discussion on calibration metrics. To keep the scale of our experiments manageable, we highlight only the latest algorithms of each sub-group used in *Peacock* (e.g., AdaFocal).

3 BACKGROUND

3.1 DEEP NEURAL NETWORK CALIBRATION

Consider a classification problem over an input feature space X and output space Y , where N labelled i.i.d pairs $(x_i, y_i)_{i=1}^N$ are randomly sampled from a training set \mathcal{D} . The model/hypothesis is then simply a mapping $h_\theta : X \rightarrow Y$, where $Y \in [0, 1]$ and θ denotes a deep neural network consisting of K neurons. The model is tasked to estimate a valid posterior such that $\sum_{k=1}^K P_i(y_k|x) = 1$, with the predicted top-1 class label $\hat{y} := \arg \max_i h_i^\theta(x)$ obtained from the logits with the top softmax confidence $\hat{P}(h^\theta) := \max_k P_i(y_k|x)$. The model is considered *perfectly calibrated* if and only if its confidence matches its probability of being correct, satisfying the formal definition $\mathbb{P}(\hat{y} = y | \hat{P} = P) = P \quad \forall P \in [0, 1]$. As this definition of calibration cannot be computed with finite samples, the most widely used approximation is the expected calibration error (ECE) (Naeini et al., 2015):

Definition 3.1 (Expected Calibration Error) *The empirical expected calibration error of a single hypothesis $h^\theta(x)$ can be written as Eq.3 in (Zhang et al., 2020) and Eq.7 in (Yang et al., 2023):*

$$\text{ECE}^d(h^\theta) = \sum_{b=1}^B \frac{n_b}{N} \|\bar{P}(h_b^\theta(x)) - \bar{y}_b\|_d^d \quad (1)$$

whereby the average predicted confidences $\bar{P}(h_b^\theta(x))$ and targets \bar{y}_b are partitioned into B bins, each containing n_b samples and $\|\cdot\|_d^d$ is the d -th power of the \mathcal{L}_d norm between the predictions and targets.

For OOD scenarios, the test distribution may diverge from the samples observed during training. Specifically, these OOD shifts can be caused by either concept shifts to the classes (changes in the posterior distribution $P(Y|X)$) or covariate shifts to input features (changes in the marginal distribution $P(X)$) (Shen et al., 2021). These OOD shifts tend to degrade model accuracy and calibration (Ovadia et al., 2019) which can be problematic for deployment. Unfortunately, achieving good calibration on both ID and OOD problem sets is non-trivial, since OOD samples typically vary

greatly from ID samples with the type and magnitude of shift unknown (Neo et al., 2024). In this work, our focus is on the problem of covariate shifts, with the goal of achieving good top-1 calibration and generalization across both ID and OOD settings.

3.2 CALIBRATION ALGORITHMS AND TECHNIQUES

Although many calibration algorithms have been proposed, each of these works tackle fundamentally different issues, with varying results and no consensus on which approach is the best. We diagnose the prevalent issues in deep neural network calibration and highlight the approaches of seven SOTA algorithms.

Adaptive Focal Parameter Selection The Focal loss (FL) (Lin et al., 2017) has been a pivotal contribution in network calibration (Mukhoti et al., 2020). As a trade-off between minimizing the Kullback-Leibler divergence and maximizing entropy, the FL: $\mathcal{L}_F = -\sum_k (1 - P_i(y_k|x_i))^\gamma \log P_i(y_k|x_i)$ is sensitive to the hyper-parameter γ , which controls the convexity of the entropy term. While strictly setting $\gamma > 1$ reduces over-confidence, it can also cause under-confidence. To circumvent this, Adaptive FL (AdaFocal) (Ghosh et al., 2022) conditionally switches between the FL and Inverse FL (Wang et al., 2021) with different selected values of γ .

$$\mathcal{L}_{\text{Ada}} = \begin{cases} -\sum_k (1 - P_i(y_k|x))^\gamma \log P_i(y_k|x) & \text{if } \gamma_{t,b} \geq 0 \\ -\sum_k (1 + P_i(y_k|x))^{|\gamma_{t,b}|} \log P_i(y_k|x) & \text{if } \gamma_{t,b} < 0, \end{cases} \quad (2)$$

Maximum Entropy Constraints Based on the Principle of Maximum Entropy (Jaynes, 1957) and an extension of the FL. MaxEnt loss (Neo et al., 2024) is designed to handle OOD samples using statistical constraints computed from the prior distribution of the training set.

$$\begin{aligned} \mathcal{L}_M^{ME} = & -\sum_k (1 - P_i(y_k|x))^\gamma \log P_i(y_k|x) \\ & + \lambda_\mu \left[\underbrace{\sum_k f(\mathcal{Y}) P_i(y_k|x) - \mu_G}_{\text{Global Expectation}} + \underbrace{\sum_k f(\mathcal{Y}) P_i(y_k|x) - \mu_{Lk}}_{\text{Local Expectation}} \right] \end{aligned} \quad (3)$$

Whereby the global expectations are computed from the entire training set such as $\mathbb{E}[\mathcal{Y}] = \sum_k P_i(y_k|x) f(\mathcal{Y}) = \mu_G$ and the local expectations are computed sample-wise from the class value characteristic function $f(\mathcal{Y})$. The Lagrange multiplier λ_μ controls the strength of the constraints, which can be solved cheaply using a numerical root-finder.

Under- and Over-confidence Trade-off A caveat to FL and its extensions alike, is that maximizing the entropy term tends to penalize all output predictions, causing under-confidence (Charoenthanakdee et al., 2021). Dual FL (Tao et al., 2023) maximizes the gap between the ground truth $P_i(y_{GT}|x)$ and the highest confidence $P_i(y_j|x)$ after the arg max class, balancing the trade-off between over- and under-confident predictions.

$$\begin{aligned} \mathcal{L}_{\text{Dual}} = & -\sum_k (1 - P_i(y_k|x) + P_i(y_j|x))^\gamma \log P_i(y_k|x) \\ \text{where } & P_i(y_j|x) = \max_k \{P_i(y_k|x) | P_i(y_k|x) < P_i(y_{GT}|x)\} \end{aligned} \quad (4)$$

Pairwise Binary Discriminatory Constraints As binary problems are easier to calibrate, CPC loss (Cheng & Vasconcelos, 2022) proposes to decompose the original multi-class problem into $\frac{K(K-1)}{2}$ binary classification problems. Whereby the predictions $P_i(y_k|x)$ are calibrated against the confidences $P_i(y_l|x)$ of the remaining $(K-1)$ pairs that do not involve the true class:

$$\mathcal{L}_{\text{CPC}}^{\text{lvl}} = -\frac{1}{(K-1)} \sum_{l \neq k} \log \frac{P_i(y_k|x)}{P_i(y_k|x) + P_i(y_l|x)} \quad (5)$$

Conditional Label Smoothing Label smoothing (LS) (Müller et al., 2019) improves calibration by artificially softening targets with a constant margin ϵ . However, LS often leads to under-confident predictions and requires time-consuming grid searches to find an optimal ϵ . To address these limitations, several approaches have proposed adaptive or conditional label smoothing functions (see Appendix A). Building upon these methods, Adaptive Conditional Label Smoothing (ACLS) (Park et al., 2023) aims to dynamically approximate the label smoothing function.

$$\mathcal{L}_{\text{ACLS}} = \begin{cases} \lambda_1 \max(0, h_k^\theta(x) - \min_k(h_k^\theta(x)) - m_{\text{ACLS}})^2 & \text{if } k = \hat{y} \\ \lambda_2 \max(0, h_{\hat{y}}^\theta(x) - h_k^\theta(x) - m_{\text{ACLS}})^2 & \text{if } k \neq \hat{y} \end{cases} \quad (6)$$

When $k = \hat{y}$, the smoothing function is directly proportional to $h_k^\theta(x)$, thereby lowering confidences. Similarly, when $k \neq \hat{y}$, the effects of the smoothing function decreases, allowing the logits and confidences to increase. m_{ACLS} denotes the ACLS margin and λ_1, λ_2 are hyperparameters for cases when $k = \hat{y}$ and $k \neq \hat{y}$.

Feature and Label Regularization Mixup (Zhang et al., 2018) is highly effective for network calibration (Thulasidasan et al., 2019; Chidambaram & Ge, 2024; Zhang et al., 2022). By interpolating a pair of inputs (x_i, x_j) and targets (y_i, y_j) , the augmented inputs and smoothed labels (\tilde{x}, \tilde{y}) are obtained using the following equations:

$$\begin{aligned} \tilde{x} &= \beta x_i + (1 - \beta)x_j \\ \tilde{y} &= \beta y_i + (1 - \beta)y_j \end{aligned} \quad (7)$$

where $\beta \in [0 - 1] \sim \text{Beta}(\alpha, \alpha)$ is a blending coefficient, randomly drawn from a Beta distribution. By considering the ordinal ranking of training samples, RankMixup (Noh et al., 2023) further improves vanilla mixup by enforcing the confidences of interpolated samples to be lower than the confidences of original samples. The ordinal relationship between “easy” and “hard” samples is maintained by a margin m_{MRL} .

$$\mathcal{L}_{\text{MRL}} = \max(0, \max_k \tilde{P}_i(\tilde{y}|\tilde{x}) - \max_k P_i(y|x) + m_{\text{MRL}}) \quad (8)$$

RankMixup (Noh et al., 2023) can be computationally inefficient due to its requirement for two forward passes: one for the original samples $P_i(y|x)$ and another for the mixed samples $\tilde{P}_i(\tilde{y}|\tilde{x})$. To improve computational efficiency, we propose an optimized version of RankMixup within *Peacock* that performs image and label mixing *batchwise*, enabling a single forward pass and faster compute times for $\tilde{P}_i(\tilde{y}|\tilde{x})$. Additional details for speeding up RankMixup can be found in Appendix B.2.

Adaptive Temperature Scaling As a post-hoc method, temperature scaling (TS) (Platt & Karampatziakis, 2007) manipulates the predictions by a scalar $\mathcal{T} \in \mathcal{R}^+$. Similar to LS, TS tends to reduce the confidence of every sample - even for correct predictions and finding a suitable \mathcal{T} requires a grid-search over a separate validation set. Adaptive Temperature scaling (AdaTS) (Joy et al., 2023) aims to learn samplewise temperatures from the features $h^\theta(x)$. By jointly learning a conditional variational autoencoder (Kingma & Welling, 2014) and a multi-layer perceptron ϕ , the samplewise temperatures are obtained as a post-processing step.

$$\mathcal{L}_{\text{AdaTS}} = -\text{ELBO}[h_i^\theta(x)] - \log\left(\frac{\exp(h_i^\theta(x)/\mathcal{T}_i)}{\sum_{k=1}^K \exp(h_k^\theta(x)/\mathcal{T}_i)}\right) \quad (9)$$

4 MOTIVATION AND PEACOCK (PUTTING IT ALL TOGETHER)

Component Synergy Why combine calibration methods? Despite their diverse approaches, all calibration algorithms discussed in Section 3.2 share the common goal of enhancing model cali-

270 bration. This suggests that they can be effectively integrated into a unified framework, leveraging
 271 their complementary strengths to achieve even better results. This section outlines our theoretical
 272 motivation for unifying calibration algorithms into *Peacock*. We first demonstrate how their equal
 273 combination improves calibration performance, subsequently we propose a novel weighted impor-
 274 tance formulation that dynamically balances loss terms to further boost performance.

276 4.1 EQUAL IMPORTANCE FORMULATION

277 Consider a multi-objective function $\mathcal{L}(\theta) = \frac{1}{A} \sum_{t=1}^A \mathcal{L}_t(\theta)$ comprising of a linear, equally weighted
 278 sum of A correlated loss terms/algorithms. From Definition 3.1, each empirical loss term $\mathcal{L}_t(\theta) \triangleq$
 279 $\frac{1}{N} \sum_i \mathcal{L}(\hat{P}_i(h_t^\theta), y_i)$ yields an individual hypothesis $\hat{P}(h_t^\theta)$ with a corresponding calibration error
 280 $\hat{P}(h_t^\theta) = \bar{y} + \text{ECE}^d(h_t^\theta)$. The unified hypothesis H^θ of the multi-objective learner $\mathcal{L}(\theta)$
 281 can then be interpreted as the average of each individual hypothesis $\bar{P}(H^\theta) = \frac{1}{A} \sum_{t=1}^A \hat{P}(h_t^\theta) =$
 282 $\bar{y} + \text{ECE}^d(H^\theta)$. When $d = 2$, the averaged squared ECE across all individual hypotheses is given
 283 by:
 284

$$285 \overline{\text{ECE}}^2(h^\theta) = \frac{1}{A} \sum_{t=1}^A \text{ECE}^2(h_t^\theta) = \frac{\text{ECE}(h_1^\theta)^2 + \text{ECE}(h_2^\theta)^2 + \dots + \text{ECE}(h_t^\theta)^2}{A} \quad (10)$$

288 As we equally consider the contributions of each individual hypotheses/loss term, with some rear-
 289 rangement the expected squared ECE of the unified multi-objective learner can be obtained as:

$$290 \text{ECE}^2(H^\theta) = \mathbb{E}_x \left[\left(\frac{1}{A} \sum_{t=1}^A \text{ECE}(h_t^\theta(x)) \right)^2 \right] = \int \left(\frac{1}{A} \sum_{t=1}^A \text{ECE}(h_t^\theta(x)) \right)^2 p(x) dx \quad (11)$$

294 where $p(x)$ is the prior probability of each input. Then from Eq. (10) and Eq. (11), the combined
 295 learner is safely bounded by the averaged squared ECE of all individual hypotheses.

$$296 \text{ECE}^2(H^\theta) \leq \overline{\text{ECE}}^2(h^\theta) \quad (12)$$

298 As the ECE^1 and ECE^2 are highly correlated (Zhang et al., 2020), we expect the upper bound in
 299 Eq. (12) to hold. This upper bound remains applicable even for temperature-scaled variants of each
 300 hypothesis, where \mathcal{T} is a temperature function.

$$301 \text{ECE}^2(\mathcal{T}(H^\theta)) \leq \overline{\text{ECE}}^2(\mathcal{T}(h^\theta)) \quad (13)$$

303 *Proof.* See Appendix C.

304 Similar to model ensembles (Zhou, 2012), loss ensembles allows a single model to perform well
 305 on multiple tasks, with additional practical benefits, such as sharing lower-level features and better
 306 compute times (Dosovitskiy & Djolonga, 2020). However, loss terms can often be conflicting,
 307 requiring trade-offs between different objectives.

309 4.2 WEIGHTED IMPORTANCE FORMULATION

310 To address these challenges, we propose a weighted importance formulation in the following section.
 311 This approach aims to find a suitable set of weights that optimizes the overall performance of the
 312 multi-objective learner. The weighted multi-objective optimization problem generally yields the
 313 following minimization problem:
 314

$$315 \min_{\theta} \mathcal{L}(\theta) = \sum_{t=1}^A w_t \mathcal{L}_t(\theta) \quad (14)$$

319 where w_t are a set of unknown scalar weights controlling each loss term. In many cases, obtaining a
 320 suitable set of weights for Eq. (14) is highly desirable. However, common approaches would either
 321 typically require expensive grid searches or predefined heuristics (Kendall et al., 2018; Chen et al.,
 322 2018). A well-studied approach in multi-objective optimization are Pareto optimal solutions, which
 323 delivers different trade-offs amongst loss terms. The goal of achieving Pareto optimality (Sener &
 Koltun, 2018; Lin et al., 2019) is defined with the following necessary conditions.

Definition 4.1 (Conditions for Pareto Optimal Calibration)

1. **Pareto dominance** A solution θ dominates another solution $\bar{\theta}$ where $\theta \prec \bar{\theta}$, if $\mathcal{L}^t(\theta) \leq \mathcal{L}^t(\bar{\theta})$ for all objectives t and $\mathbf{L}(\theta_1, \dots, \theta_A) \neq \mathbf{L}(\bar{\theta}_1, \dots, \bar{\theta}_A)$.
2. **Pareto optimality** Solution θ^* is considered Pareto optimal if there exists no other solution θ that dominates θ^* such that $\theta \prec \theta^*$.

Assuming loss terms are convex and optimizable with gradient descent, Pareto optimal weights for each loss can be obtained through the Karush-Kuhn-Tucker (KKT) conditions (Fliege & Svaiter, 2000; Schäffler et al., 2002), by minimizing the following objective (Désidéri, 2012; Sener & Koltun, 2018):

$$\min_{w_1, \dots, w_t} \left\{ \left\| \sum_{t=1}^A w_t \nabla_{\theta} \mathcal{L}_t(\theta) \right\|_2 \left| \sum_{t=1}^A w_t = 1, w_t \geq 0 \quad \forall t \right. \right\} \quad (15)$$

Previous works (Désidéri, 2012; Sener & Koltun, 2018) have shown that the solution to Eq. (15) is either zero or provides a gradient direction that improves all loss components.

Practical Considerations Generally, optimizing for w_t in Eq. (15) requires a separate optimizer and the recomputation of the gradients $\nabla_{\theta} \mathcal{L}_t(\theta)$ for each loss term. This involves retaining the computational graph³ for A backward passes, which slows down training speeds and grows prohibitively more expensive as A becomes larger.

To circumvent this, we propose a fast, elegant and efficient alternative to recomputing gradients, by replacing $\nabla_{\theta} \mathcal{L}_t(\theta)$ with decrease rate estimates for each loss term. Assuming model parameters θ are updated via $\theta' \leftarrow \theta - \eta \nabla_{\theta} \mathcal{L}_t(\theta)$, we propose to balance learning loss terms with the following direction-orientated objective:

$$\min_{w_1, \dots, w_t} \left\{ \left\| \sum_{t=1}^A w_t \sqrt{\frac{\Delta_{\theta} \mathcal{L}_t(\theta)}{\eta}} \right\|_2 \left| \sum_{t=1}^A w_t = 1, w_t \geq 0 \quad \forall t \right. \right\} \quad (16)$$

Proof. See Appendix C.2.

where the decrease rate estimates $\sqrt{\frac{\Delta_{\theta} \mathcal{L}_t(\theta)}{\eta}}$, are derived using the first-order Taylor approximation with a sufficiently small step size η . As long as the KKT conditions are satisfied and loss terms $\mathcal{L}_t(\theta)$ are monotonically decreasing, faster performance can be achieved using Eq. (16). With the added benefit of $w_t \propto \sqrt{\frac{\Delta_{\theta} \mathcal{L}_t(\theta)}{\eta}}$ which ensures balanced learning rates across all loss terms, preventing any single term from dominating the optimization process.

Direction Weighted Self-Attention To optimize the objective in Eq. (16), we propose using a direction weighted self-attention block. Fig. 3 illustrates our self-attention block, which accepts an array of loss terms $\mathcal{L}_t(\theta)$ as inputs and outputs a set of weights w_t . The Value (V), Key (K) and Query (Q) neurons are of size $A \times A$ and the softmax function σ is applied to ensure that $\sum_t w_t = 1$. The learning dynamics of the direction weighted self-attention block are discussed in Appendix B.3.

Full details can be found in Algorithm 1, which shows all calibration components of *Peacock* and an optional step for obtaining importance weights. A peculiar finding in our ablation study, is that removing ACLS tends to lead to better performance in *Peacock*. Since $\mathcal{L}_{\text{AdaFocal}}^{\text{Dual}}$ contains the CE loss, we only apply importance weights to the auxiliary loss terms with the final objective function given by: $\mathcal{L}^{\text{Peacock}} = \mathcal{L}_{\text{AdaFocal}}^{\text{Dual}} + w_1 \mathcal{L}_{\text{constraints}}^{\text{ME}} + w_2 \mathcal{L}_{\text{CPC}}^{\text{lv1}} + w_3 \mathcal{L}_{\text{MRL}}$.

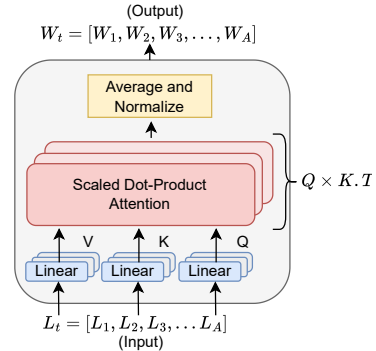


Figure 3: Our direction weighted self-attention block learns the importance of each loss term.

³ For more details, see Pytorch autograd framework: <https://pytorch.org/docs/stable/autograd.html>

Dataset	Metric	CE	MaxEnt	AdaFocal	RankMixup	CPC	Dual	ACLS	Peacock (Eq.)	Peacock (Impt.)
CIFAR10-C	Acc. ↑	77.9±0.3	78.3±0.2	77.7±0.3	77.8±0.3	77.6±0.2	77.9±0.3	78.1±0.4	76.8±0.2	77.3±0.4
	ECE ↓	14.5±0.4	6.5±0.2	6.8±0.2	11.9±0.3	11.2±0.2	7.4±0.3	11.2±0.4	<u>6.3</u> ±0.1	6.2 ±0.3
	CECE ↓	3.3±0.1	2.6±0.1	2.6±0.1	2.9±0.1	2.8±0.1	2.7±0.1	2.8±0.1	<u>2.7</u> ±0.1	2.6 ±0.1
	KSE ↓	14.5±0.4	6.2±0.1	6.5±0.3	11.8±0.2	11.1±0.1	7.0±0.3	11.3±0.4	<u>6.1</u> ±0.2	6.1 ±0.3
CIFAR100-C	Acc. ↑	52.5±0.1	52.4±0.1	52.9±0.2	52.3±0.1	52.0±0.1	52.8±0.1	52.6±0.1	51.8±0.1	52.6±0.1
	ECE ↓	10.6±0.1	11.6±0.6	13.6±0.1	11.0±1.4	13.2±0.5	15.7±0.1	12.2±0.2	<u>9.6</u> ±0.2	9.3 ±0.1
	CECE ↓	0.4±0.1	0.5±0.1	0.5±0.1	0.4±0.1	0.4±0.1	0.5±0.1	0.4±0.1	<u>0.4</u> ±0.1	0.4 ±0.1
	KSE ↓	9.3±0.1	11.4±0.6	13.3±0.1	10.3±1.5	11.5±0.5	15.3±0.1	11.7±0.3	<u>9.6</u> ±0.4	9.2 ±0.6
TinyImageNet-C	Acc. ↑	25.2±0.1	22.0±0.1	25.1±0.1	23.1±0.3	23.7±0.1	22.9±0.6	22.1±0.1	23.3±0.5	23.6±0.2
	ECE ↓	15.7±0.5	12.8±0.1	13.8±0.4	20.2±0.1	16.0±0.5	19.2±0.4	19.8±0.2	<u>10.6</u> ±0.2	10.4 ±0.2
	CECE ↓	0.3±0.1	0.3±0.1	0.3±0.1	0.4±0.1	0.3±0.1	0.3±0.1	0.3±0.1	<u>0.3</u> ±0.1	0.3 ±0.1
	KSE ↓	15.7±0.5	12.8±0.2	13.8±0.3	20.2±0.2	15.7±0.2	19.2±0.7	19.8±0.1	<u>10.6</u> ±0.2	10.3 ±0.2
Camelyon17	Acc. ↑	81.7±0.7	79.7±1.7	74.5±0.1	74.9±4.0	77.7±1.6	78.4±2.9	77.0±1.2	79.3±2.7	83.2±1.1
	ECE ↓	15.5±1.1	12.4±0.1	20.4±0.4	22.4±4.6	20.2±1.6	15.4±2.5	19.8±0.2	<u>11.7</u> ±0.7	9.8 ±1.8
	CECE ↓	16.7±1.3	16.2±0.1	23.7±0.7	23.6±4.8	21.3±2.0	20.3±2.9	22.0±1.1	<u>14.0</u> ±0.4	13.7 ±1.6
	KSE ↓	15.5±1.1	12.3±0.8	20.4±0.1	22.4±4.6	20.2±1.6	15.4±2.4	19.6±0.1	<u>11.7</u> ±0.7	9.8 ±1.8
iWildCam	Acc. ↑	52.2±0.3	50.9±1.0	54.1±1.7	56.7±0.3	55.8±2.2	54.6±2.1	55.2±2.2	51.7±0.8	54.5±0.8
	ECE ↓	30.6±0.8	21.0±3.2	23.0±0.5	25.5±0.7	20.3±1.1	13.0±2.5	20.6±1.8	9.7 ±0.3	<u>12.6</u> ±1.4
	CECE ↓	0.4±0.1	0.4±0.1	0.3±0.1	0.4±0.1	0.3±0.1	0.4±0.1	0.3±0.1	0.3 ±0.1	<u>0.3</u> ±0.1
	KSE ↓	30.6±0.8	21.0±3.2	23.0±0.5	25.5±0.7	19.5±1.6	13.0±2.5	20.6±1.8	9.7 ±0.3	<u>12.6</u> ±1.4
FmoW	Acc. ↑	35.1±0.5	33.5±0.1	35.5±0.7	35.8±0.1	36.4±0.1	35.1±0.2	37.5±0.1	35.1±0.2	35.5±0.1
	ECE ↓	39.8±0.2	20.0±9.9	20.9±8.6	41.7±0.1	22.4±0.9	10.7±0.1	21.7±0.2	<u>10.6</u> ±0.4	10.5 ±0.3
	CECE ↓	1.5±0.1	1.0±0.3	1.0±0.2	1.5±0.1	0.9±0.1	0.6±0.1	0.9±0.1	<u>0.6</u> ±0.1	0.6 ±0.1
	KSE ↓	39.8±0.2	20.0±9.9	20.9±8.6	41.7±0.1	22.4±0.9	10.7±0.1	21.7±0.2	<u>10.6</u> ±0.4	10.5 ±0.3
Amazon	Acc. ↑	55.8±0.3	64.6±0.4	59.6±2.7	56.9±0.1	56.9±0.1	60.7±3.8	56.9±0.2	57.5±0.6	64.9±0.8
	ECE ↓	7.0±0.5	5.0±0.6	6.7±1.0	43.1±0.1	7.4±0.5	5.8±2.3	42.0±0.1	6.5±3.2	5.0 ±1.1
	CECE ↓	6.4±0.3	3.8±0.6	3.3±0.8	17.2±0.1	4.8±0.2	2.5±0.9	16.8±0.1	<u>2.9</u> ±1.3	2.3 ±0.3
	KSE ↓	7.0±0.5	5.1±0.6	7.5±0.6	43.1±0.3	10.9±3.3	8.1±0.1	42.1±0.1	8.6±1.0	<u>6.6</u> ±0.4
CivilComments	Acc. ↑	90.3±1.0	91.3±0.1	91.4±0.1	88.6±1.0	88.6±0.8	91.5±0.1	88.6±0.1	90.1±0.7	90.8±0.5
	ECE ↓	10.4±0.4	4.8±0.2	7.8±1.7	11.4±0.5	2.4±0.4	4.2±0.1	11.1±0.1	2.1 ±0.8	<u>4.2</u> ±1.0
	CECE ↓	10.4±0.5	5.6±0.1	8.1±1.8	11.4±0.5	2.4±0.4	4.8±0.3	11.1±0.2	2.2 ±0.3	<u>4.6</u> ±0.8
	KSE ↓	10.4±0.4	6.7±0.1	7.7±1.7	5.8±0.1	2.9±0.4	4.3±0.1	11.0±0.1	3.3 ±0.8	<u>4.2</u> ±1.0

Table 1: We report the OOD test scores (%) computed across 3 seeds, evaluated on both synthetic and wild benchmarks for *Peacock* and recent baselines. *Peacock* greatly improves calibration and maintains model accuracy. The best calibration scores in bold, second best are underlined.

5 EXPERIMENTS AND ANALYSIS

5.1 EXPERIMENT SETUP

Evaluation Metrics Following (Guo et al., 2017; Mukhoti et al., 2020; Neo et al., 2024), we use the Expected Calibration Error (ECE), Classwise Calibration Error (CECE) (Nixon et al., 2019) and Kolmogorov-Smirnov Error (KSE) (Gupta et al., 2021) for evaluation. For fair comparisons, we follow the evaluation protocols of other authors and compute calibration errors using 15 bins with the mean and standard deviation shown across seeds. Additional details of each metric are included in Appendix A.5.

Datasets We evaluate *Peacock* on a total of eight OOD image and text benchmarks. For synthetic datasets, we use CIFAR (Krishnan & Tickoo, 2020) and TinyImageNet (Deng et al., 2009) for training/validation and CIFAR-C/TinyImageNet-C (Hendrycks & Dietterich, 2019) for testing. For Wild datasets, we use Camelyon-17 (Bandi et al., 2019), iWildCam (Beery et al., 2020), FmoW (Christie et al., 2018), Amazon (Ni et al., 2019) and CivilComments (Borkan et al., 2019) from the Wilds benchmark (Koh et al., 2021). OOD data is never used for training or validating a model, only for testing.

Baselines We compare equal and importance weighted *Peacock* against an uncalibrated baseline (CE) and six components, specifically MaxEnt (Neo et al., 2024), AdaFocal (Ghosh et al., 2022), RankMixup (Noh et al., 2023), CPC (Cheng & Vasconcelos, 2022), Dual (Tao et al., 2023), ACLS (Park et al., 2023). For our analysis on image tasks, we use ResNet-18, ResNet-50 (He et al., 2016), SWINV2 (Liu et al., 2022b) and RoBERTa (Liu et al., 2019b) for text tasks. We perform post-hoc processing with AdaTS (Joy et al., 2023) and compare different weighted formulations analyzing the overall contributions of each component used in *Peacock*. For additional details of each dataset task, hyper-parameters and illustrations of synthetic and wild OOD shifts, we refer readers to Appendix D.1.

Dataset	ECE ↓	CE	MaxEnt	AdaFocal	RankMixup	CPC	Dual	ACLS	Peacock (Eq.)	Peacock (Impt.)
CIFAR10-C	Pre	14.5±0.4	6.5±0.2	6.8±0.2	11.9±0.3	11.2±0.2	7.4±0.3	11.2±0.4	6.3±0.1	6.2±0.3
	Post	7.5±0.1	6.9±0.2	6.9±0.4	7.1±0.3	7.0±0.1	7.3±0.3	7.3±0.1	6.6±0.1	6.9±0.3
	Avg.	11.0±0.3	6.7±0.2	6.9±0.3	9.5±0.3	9.1±0.2	7.4±0.3	9.3±0.3	6.4±0.1	6.6±0.3
CIFAR100-C	Pre	10.6±0.1	11.6±0.6	13.6±0.1	11.0±1.4	13.2±0.5	15.7±0.1	12.2±0.2	9.6±0.2	9.3±0.1
	Post	8.1±0.4	7.2±0.1	7.5±0.1	8.0±0.1	10.2±0.2	8.4±0.2	8.1±0.4	8.7±0.2	8.9±0.2
	Avg.	9.4±0.3	9.4±0.4	10.6±0.1	9.5±1.0	11.7±0.3	12.1±0.2	10.2±0.3	9.2±0.2	9.1±0.2
TinyImageNet-C	Pre	15.7±0.5	12.8±0.1	13.8±0.4	20.2±0.1	16.0±0.5	19.2±0.4	19.8±0.2	10.6±0.2	10.4±0.2
	Post	25.4±0.3	20.2±0.2	26.1±0.4	24.3±0.3	20.7±0.4	24.3±0.3	24.4±0.4	11.9±0.2	12.6±0.2
	Avg.	20.5±0.4	16.5±0.2	20.0±0.4	22.3±0.3	18.4±0.5	21.8±0.4	22.1±0.3	11.2±0.4	11.5±0.4
Camelyon17	Pre	15.5±1.1	12.4±0.1	13.8±0.4	20.2±0.1	16.0±0.5	19.2±0.4	19.8±0.2	11.7±0.7	9.87±1.8
	Post	33.1±0.4	11.2±1.1	15.4±1.5	14.9±1.4	12.2±1.4	18.1±0.3	11.6±0.7	9.87±3.0	12.2±1.7
	Avg.	24.3±0.8	11.8±0.6	14.6±1.0	17.6±0.8	14.1±0.9	18.7±0.4	15.7±0.4	10.8±1.5	11.0±1.8
iWildCam	Pre	30.6±0.8	21.0±3.2	23.0±0.5	25.5±0.7	20.3±1.1	13.0±2.5	20.6±1.8	12.0±0.1	11.7±2.3
	Post	8.0±2.2	8.6±0.6	7.9±1.3	9.5±2.5	11.7±1.9	8.0±1.0	7.9±0.4	8.9±0.6	6.7±0.5
	Avg.	19.3±1.5	14.8±2.4	15.5±1.1	17.5±1.6	16.0±1.6	10.5±1.6	14.3±1.2	10.5±0.4	9.2±1.4
FmoW	Pre	39.8±0.2	20.0±9.9	20.9±8.6	41.7±0.1	22.4±0.9	10.7±0.1	21.7±0.2	10.6±0.4	10.5±0.3
	Post	25.6±0.6	4.9±0.9	6.2±0.5	7.9±0.3	7.6±0.9	5.6±0.9	6.1±0.5	5.3±0.8	5.9±0.5
	Avg.	32.7±0.5	12.5±5.5	13.6±4.7	24.8±0.5	15.0±0.8	8.2±0.2	13.9±0.2	7.9±0.4	8.2±0.8

Table 2: ECE (%) scores before and after AdaTS (Joy et al., 2023) for the different OOD datasets. *Peacock* delivers the best overall calibration performance, despite using temperatures obtained ID.

5.2 COMPARISONS TO PUBLISHED BASELINES

In-Distribution Performance Fig. 1 and Table 5 showcase the ID results, demonstrating *Peacock*’s consistently strong calibration performance across datasets. Following Eq. (12), *Peacock*’s ECE is significantly lower than the empirical average $\overline{\text{ECE}}$ of all methods. While individual algorithms may vary in performance across datasets, we demonstrate that combining them through *Peacock* consistently improves calibration - with no significant loss in accuracy. Additional ID results and discussions are included in Appendix D.2.

Out-of-Distribution Performance Similar to the ID results, each individual algorithm’s independent performance varies across different datasets. Table 1 shows that by combining algorithms through equal importance *Peacock*, consistent improvements in OOD performance can be achieved on both synthetic and real-world image and text datasets. Our findings further demonstrate that the equal importance formulation of *Peacock* also adheres to the theoretical upper bound in Eq. (12) for OOD performance. Finally, *Peacock* can be further enhanced using our weighted importance formulation, achieving improved results and good generalization properties across all datasets.

Training Time per Epoch Fig. 4 shows the average wall-clock time per epoch (forward pass, loss-calculation and back-propagation) for each method trained on CIFAR10. We report the wall-clock time in seconds on a NVIDIA GeForce RTX 2070 GPU with i7-10700 CPU. In general, each algorithm has similar speeds, with RankMixup taking the longest since another forward pass is required to obtain the logits of interpolated samples. On the other hand, *Peacock* is optimized (see Appendix B.2) to remain competitive with other baselines, despite combining multiple algorithms together.

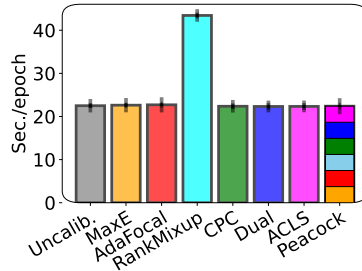


Figure 4: Wall-clock time for each method on CIFAR10. *Peacock* is as fast as each of its components.

5.3 POST-HOC PROCESSING

For post-hoc calibration, we apply AdaTS (Joy et al., 2023) to each method. The samplewise temperatures are obtained from an ID validation set and applied to the OOD test sets. Table 2 presents the ECE scores of each algorithm before and after applying AdaTS for all six OOD image datasets. Our findings demonstrate that *Peacock* delivers the best overall calibration performance, both before and after temperature scaling. In cases where *Peacock* does not deliver the best calibration, we can see that its performance is relatively close to the best score. While AdaTS generally improves OOD ECE, applying it to already well-calibrated models can sometimes lead to degraded performance. For instance, MaxEnt Loss achieves the best OOD calibration with 4.9% on FmoW without AdaTS, but applying it subsequently would cause the ECE to worsen to 12.5%. This discrepancy can be

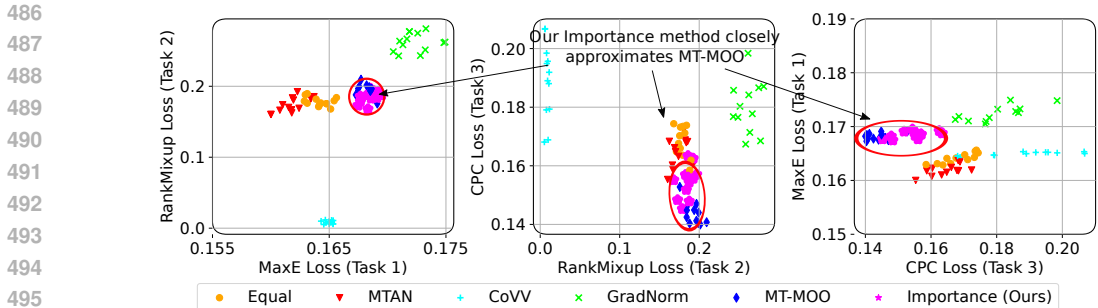


Figure 5: Solution given by different MOO algorithms for each of the auxiliary test losses. Our importance formulation effectively balances trade-offs between loss terms. Bottom-left is better.

Algorithm	Acc (%)	ECE (%)	Speed (Sec)	w1	w2	w3	$\sum_t^A w_t = 1$
Equal-Importance	51.8 \pm 0.1	9.6 \pm 0.2	50.1 \pm 0.2	0.33	0.33	0.33	Yes
MTAN (Liu et al., 2019a)	50.8 \pm 0.1	10.2 \pm 0.4	50.7 \pm 0.2	0.33	0.33	0.33	Yes
CoVV (Groenendijk et al., 2021)	52.6 \pm 0.1	12.0 \pm 0.4	50.6 \pm 0.2	0.01	0.52	0.47	Yes
GradNorm (Chen et al., 2018)	48.9 \pm 0.6	9.3 \pm 0.7	85.8 \pm 0.7	2.99	0.00	0.00	No
MT-MOO (Sener & Koltun, 2018)	51.9 \pm 0.5	9.3 \pm 0.5	67.5 \pm 0.5	0.36	0.47	0.16	Yes
Weighted-Importance (Ours)	52.6 \pm 0.1	9.3 \pm 0.3	50.5 \pm 0.2	0.00	0.51	0.49	Yes

Table 3: Comparisons of different multi-objective optimization methods for *Peacock*. Our weighted importance formulation is fast and effective.

attributed to the disconnect between the training, validation sets and test set, explaining the higher calibration errors after temperature scaling (Ovadia et al., 2019). However, by combining multiple calibration algorithms together, *Peacock* displays the best generalization behavior even when using temperatures obtained ID. As AdaTS is designed solely for image tasks, we apply vanilla TS for Amazon and CivilComments indicating their results and ideal temperatures obtained from grid-search in Table 7. We further highlight that temperature scaling only manipulates the predicted confidences and does not affect recognition accuracy.

5.4 WEIGHTED PEACOCK PERFORMANCE AND ANALYSIS

We compare various weighted objective optimization methods for *Peacock*. Namely, using equal weights, MTAN (Liu et al., 2019a), GradNorm (Chen et al., 2018), CoVV (Groenendijk et al., 2021), MT-MOO (Sener & Koltun, 2018) and our proposed weighted importance variant, on CIFAR100/100-C using ResNet-18. All MOO methods are initialized with uniform weights, with the final test accuracy, ECE and weights shown upon convergence with the average training wall clock time per epoch in Table 3. Our results indicate that while each algorithm yields distinct solutions/weights for each loss term, they achieve comparable accuracies and calibration errors. GradNorm and MT-MOO have longer training times (about 65% and 35% respectively) since they require the recomputation of $\nabla_{\theta} \mathcal{L}_t(\theta)$ for each loss term. Conversely, MTAN and CoVV offers faster performance, but has higher ECE. Fig. 5 further demonstrates that the auxiliary test losses for each MOO method yield similar solutions, with MT-MOO achieving the optimal solution. Our method closely approximates MT-MOO but is faster and more efficient. We further discuss contributions of each calibration loss term in Appendix D.3 and the limitations of our work in Appendix E.

6 CONCLUSIONS

We present *Peacock*, a unified framework for neural network calibration. By formulating unification as a multi-objective optimization problem, we demonstrate that combining calibration components improves performance on both ID and OOD tasks. Our proposed weighted importance form of *Peacock* is fast and effective in delivering good Pareto Optimal performance. Despite incorporating multiple algorithms, *Peacock*’s complements post-hoc processing and remains fast in terms of computational speed. Our method shows clear performance gains with RankMixup and MaxEnt loss offering the most improvements.

REFERENCES

- 540
541
542 Peter Bandi, Oscar Geessink, Quirine Manson, Marcory Van Dijk, Maschenka Balkenhol, Meyke
543 Hermesen, Babak Ehteshami Bejnordi, Byungjae Lee, Kyunghyun Paeng, Aoxiao Zhong, et al.
544 From detection of individual metastases to classification of lymph node status at the patient level:
545 The CAMELYON17 challenge. *IEEE Transactions on Medical Imaging*, 38(2):550–560, 2019.
- 546 Sara Beery, Elijah Cole, and Arvi Gjoka. The iWildCam 2020 competition dataset. *arXiv*,
547 abs/2004.10340, 2020.
- 548 Daniel Borkan, Lucas Dixon, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. Nuanced
549 metrics for measuring unintended bias with real data for text classification. In *Compan-*
550 *ion Proceedings of The 2019 World Wide Web Conference*, WWW ’19, pp. 491–500, New
551 York, NY, USA, 2019. Association for Computing Machinery. ISBN 9781450366755. doi:
552 10.1145/3308560.3317593. URL <https://doi.org/10.1145/3308560.3317593>.
- 553 Rich Caruana. Multitask learning. *Machine Learning*, 28:41–75, 1997. URL [https://api.](https://api.semanticscholar.org/CorpusID:45998148)
554 [semanticscholar.org/CorpusID:45998148](https://api.semanticscholar.org/CorpusID:45998148).
- 555 Nontawat Charoenphakdee, Jayakorn Vongkulbhisal, Nuttapon Chairatanakul, and Masashi
556 Sugiyama. On focal loss for class-posterior probability estimation: A theoretical perspective. In
557 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,
558 pp. 5202–5211, June 2021.
- 559 Zhao Chen, Vijay Badrinarayanan, Chen-Yu Lee, and Andrew Rabinovich. GradNorm: Gradient
560 normalization for adaptive loss balancing in deep multitask networks. In Jennifer Dy and Andreas
561 Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80
562 of *Proceedings of Machine Learning Research*, pp. 794–803. PMLR, 10–15 Jul 2018. URL
563 <https://proceedings.mlr.press/v80/chen18a.html>.
- 564 Jiacheng Cheng and Nuno Vasconcelos. Calibrating deep neural networks by pairwise constraints. In
565 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,
566 pp. 13709–13718, June 2022.
- 567 Muthu Chidambaram and Rong Ge. On the limitations of temperature scaling for distributions with
568 overlaps. In *The Twelfth International Conference on Learning Representations*, 2024. URL
569 <https://openreview.net/forum?id=zavLQJlXjB>.
- 570 Gordon Christie, Neil Fendley, James Wilson, and Ryan Mukherjee. Functional map of the world.
571 In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 6172–
572 6180, 2018. doi: 10.1109/CVPR.2018.00646.
- 573 Nina Corvelo Benz and Manuel Rodriguez. Human-aligned calibration for ai-assisted decision mak-
574 ing. In A. Oh, T. Neumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in*
575 *Neural Information Processing Systems*, volume 36, pp. 14609–14636. Curran Associates, Inc.,
576 2023. URL [https://proceedings.neurips.cc/paper_files/paper/2023/](https://proceedings.neurips.cc/paper_files/paper/2023/file/2f1d1196426ba84f47d115cac3dcb9d8-Paper-Conference.pdf)
577 [file/2f1d1196426ba84f47d115cac3dcb9d8-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/2f1d1196426ba84f47d115cac3dcb9d8-Paper-Conference.pdf).
- 578 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
579 hierarchical image database. In *Proceedings of the IEEE Conference on Computer Vision and*
580 *Pattern Recognition (CVPR)*, pp. 248–255, 2009. doi: 10.1109/CVPR.2009.5206848.
- 581 Alexey Dosovitskiy and Josip Djolonga. You only train once: Loss-conditional training of deep net-
582 works. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa,*
583 *Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. URL [https://openreview.net/](https://openreview.net/forum?id=HyxY6JHKwr)
584 [forum?id=HyxY6JHKwr](https://openreview.net/forum?id=HyxY6JHKwr).
- 585 Jean-Antoine Désidéri. Multiple-gradient descent algorithm (mgda) for multiobjective optimization.
586 *Comptes Rendus Mathématique*, 350:313–318, 03 2012. doi: 10.1016/j.crma.2012.03.014.
- 587 Andre Esteva, Brett Kuprel, Roberto A. Novoa, Justin M. Ko, Susan M. Swetter, Helen M. Blau,
588 and Sebastian Thrun. Dermatologist-level classification of skin cancer with deep neural networks.
589 *Nature*, 542:115–118, 2017.
- 590
591
592
593

- 594 Jörg Fliege and Benar Fux Svaiter. Steepest descent methods for multicriteria optimization.
595 *Mathematical Methods of Operations Research*, 51:479–494, 2000. URL [https://api.
596 semanticscholar.org/CorpusID:44256411](https://api.semanticscholar.org/CorpusID:44256411).
- 597 Arindam Ghosh, Thomas Schaaf, and Matthew Gormley. Adafocal: Calibration-aware adaptive fo-
598 cal loss. In *Advances in Neural Information Processing Systems*, volume 35, pp. 1583–1595,
599 2022. URL [https://proceedings.neurips.cc/paper_files/paper/2022/
600 file/0a692a24dbc744fca340b9ba33bc6522-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/0a692a24dbc744fca340b9ba33bc6522-Paper-Conference.pdf).
- 601 Ting Gong, Tyler Lee, Cory Stephenson, Venkata Renduchintala, Suchismita Padhy, Anthony Ndi-
602 rango, Gokce Keskin, and Oguz Elibol. A comparison of loss weighting strategies for multi task
603 learning in deep neural networks. *IEEE Access*, PP:1–1, 09 2019. doi: 10.1109/ACCESS.2019.
604 2943604.
- 605 Rick Groenendijk, Sezer Karaoglu, Theo Gevers, and Thomas Mensink. Multi-loss weighting with
606 coefficient of variations. In *Proceedings of the IEEE/CVF Winter Conference on Applications of
607 Computer Vision*, pp. 1469–1478, 2021.
- 608 Sebastian Gruber and Florian Buettner. Better uncertainty calibration via proper
609 scores for classification and beyond. In S. Koyejo, S. Mohamed, A. Agarwal,
610 D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Process-
611 ing Systems*, volume 35, pp. 8618–8632. Curran Associates, Inc., 2022. URL
612 [https://proceedings.neurips.cc/paper_files/paper/2022/file/
613 3915a87ddac8e8c2f23dbabbcee6eec9-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/3915a87ddac8e8c2f23dbabbcee6eec9-Paper-Conference.pdf).
- 614 Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural
615 networks. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International
616 Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp.
617 1321–1330. PMLR, 06–11 Aug 2017. URL [https://proceedings.mlr.press/v70/
618 guol7a.html](https://proceedings.mlr.press/v70/guol7a.html).
- 619 Michelle Guo, Albert Haque, De-An Huang, Serena Yeung, and Li Fei-Fei. Dynamic task priori-
620 tization for multitask learning. In *Proceedings of the European Conference on Computer Vision
621 (ECCV)*, September 2018.
- 622 Kartik Gupta, Amir Rahimi, Thalaiyasingam Ajanthan, Thomas Mensink, Cristian Sminchisescu,
623 and Richard Hartley. Calibration of neural networks using splines. In *9th International Confer-
624 ence on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenRe-
625 view.net, 2021. URL <https://openreview.net/forum?id=eQe8DEWNN2W>.
- 626 Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning*.
627 Springer Series in Statistics. Springer, New York, NY, USA, 2001.
- 628 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
629 nition. In *Proc. IEEE Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, June
630 2016.
- 631 Ramya Hebbalaguppe, Jatin Prakash, Neelabh Madan, and Chetan Arora. A stitch in time saves
632 nine: A train-time regularizing loss for improved neural network calibration. In *Proceedings
633 of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 16081–
634 16090, June 2022.
- 635 Dan Hendrycks and Thomas G. Dietterich. Benchmarking neural network robustness to common
636 corruptions and perturbations. In *7th International Conference on Learning Representations,
637 ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019. URL <https://openreview.net/forum?id=HJz6tiCqYm>.
- 638 E. T. Jaynes. Information theory and statistical mechanics. *Physical Review*, 106:620–630, May
639 1957. doi: 10.1103/PhysRev.106.620.
- 640 Heinrich Jiang, Been Kim, Melody Guan, and Maya Gupta. To trust or not to trust a classi-
641 fier. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett
642 (eds.), *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc.,
643 644 645 646 647

- 648 2018. URL [https://proceedings.neurips.cc/paper_files/paper/2018/](https://proceedings.neurips.cc/paper_files/paper/2018/file/7180cffd6a8e829dacfc2a31b3f72ece-Paper.pdf)
649 [file/7180cffd6a8e829dacfc2a31b3f72ece-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2018/file/7180cffd6a8e829dacfc2a31b3f72ece-Paper.pdf).
- 650
- 651 Tom Joy, Francesco Pinto, Ser-Nam Lim, Philip H.S. Torr, and Puneet K. Dokania. Sample-
652 dependent adaptive temperature scaling for improved calibration. *Proceedings of the AAAI Con-*
653 *ference on Artificial Intelligence*, 37(12):14919–14926, Jun. 2023. doi: 10.1609/aaai.v37i12.
654 26742. URL <https://ojs.aaai.org/index.php/AAAI/article/view/26742>.
- 655 Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer
656 vision? In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan,
657 and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Cur-
658 ran Associates, Inc., 2017. URL [https://proceedings.neurips.cc/paper_files/](https://proceedings.neurips.cc/paper_files/paper/2017/file/2650d6089a6d640c5e85b2b88265dc2b-Paper.pdf)
659 [paper/2017/file/2650d6089a6d640c5e85b2b88265dc2b-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/2650d6089a6d640c5e85b2b88265dc2b-Paper.pdf).
- 660
- 661 Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weigh losses
662 for scene geometry and semantics. In *Proceedings of the IEEE Conference on Computer Vision*
663 *and Pattern Recognition (CVPR)*, June 2018.
- 664 Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In Yoshua Bengio and Yann
665 LeCun (eds.), *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB,*
666 *Canada, April 14-16, 2014, Conference Track Proceedings*, 2014. URL [http://arxiv.org/](http://arxiv.org/abs/1312.6114)
667 [abs/1312.6114](http://arxiv.org/abs/1312.6114).
- 668 Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Bal-
669 subramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, Tony Lee,
670 Etienne David, Ian Stavness, Wei Guo, Berton Earnshaw, Imran Haque, Sara M Beery, Jure
671 Leskovec, Anshul Kundaje, Emma Pierson, Sergey Levine, Chelsea Finn, and Percy Liang.
672 Wilds: A benchmark of in-the-wild distribution shifts. In Marina Meila and Tong Zhang
673 (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of
674 *Proceedings of Machine Learning Research*, pp. 5637–5664. PMLR, 18–24 Jul 2021. URL
675 <https://proceedings.mlr.press/v139/koh21a.html>.
- 676
- 677 Ranganath Krishnan and Omesh Tickoo. Improving model calibration with accuracy versus un-
678 certainty optimization. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin
679 (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 18237–18248. Cur-
680 ran Associates, Inc., 2020. URL [https://proceedings.neurips.cc/paper_files/](https://proceedings.neurips.cc/paper_files/paper/2020/file/d3d9446802a44259755d38e6d163e820-Paper.pdf)
681 [paper/2020/file/d3d9446802a44259755d38e6d163e820-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/d3d9446802a44259755d38e6d163e820-Paper.pdf).
- 682
- 683 A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. Master’s
684 thesis, Department of Computer Science, University of Toronto, 2009.
- 685
- 686 Volodymyr Kuleshov and Shachi Deshpande. Calibrated and sharp uncertainties in deep learning
687 via density estimation. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari,
688 Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International Conference on Machine*
689 *Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 11683–11693. PMLR,
690 17–23 Jul 2022. URL [https://proceedings.mlr.press/v162/kuleshov22a.](https://proceedings.mlr.press/v162/kuleshov22a.html)
691 [html](https://proceedings.mlr.press/v162/kuleshov22a.html).
- 692
- 693 Zhaoqi Leng, Mingxing Tan, Chenxi Liu, Ekin Dogus Cubuk, Jay Shi, Shuyang Cheng, and
694 Dragomir Anguelov. Polyloss: A polynomial expansion perspective of classification loss func-
695 tions. In *The Tenth International Conference on Learning Representations, ICLR 2022, Vir-*
696 *tual Event, April 25-29, 2022*. OpenReview.net, 2022. URL [https://openreview.net/](https://openreview.net/forum?id=gSdSJoenuPI)
697 [forum?id=gSdSJoenuPI](https://openreview.net/forum?id=gSdSJoenuPI).
- 698
- 699 Changsheng Li, Junchi Yan, Fan Wei, Weishan Dong, Qingshan Liu, and Hongyuan Zha. Self-
700 paced multi-task learning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 31(1),
701 Feb. 2017. doi: 10.1609/aaai.v31i1.10847. URL [https://ojs.aaai.org/index.php/](https://ojs.aaai.org/index.php/AAAI/article/view/10847)
[AAAI/article/view/10847](https://ojs.aaai.org/index.php/AAAI/article/view/10847).
- 702
- 703 Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollar. Focal loss for dense object
704 detection. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, Oct
705 2017.

- 702 Xi Lin, Hui-Ling Zhen, Zhenhua Li, Qing-Fu Zhang, and Sam Kwong. Pareto multi-task learn-
703 ing. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett
704 (eds.), *Advances in Neural Information Processing Systems*, volume 32, pp. 12037–12047. Cur-
705 ran Associates, Inc., 2019. URL [https://proceedings.neurips.cc/paper_files/
706 paper/2019/file/685bfde03eb646c27ed565881917c71c-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/685bfde03eb646c27ed565881917c71c-Paper.pdf).
- 707
- 708 Bingyuan Liu, Ismail Ben Ayed, Adrian Galdran, and Jose Dolz. The devil is in the margin: Margin-
709 based label smoothing for network calibration. In *Proceedings of the IEEE/CVF Conference on
710 Computer Vision and Pattern Recognition (CVPR)*, pp. 80–88, June 2022a.
- 711
- 712 Bingyuan Liu, Jérôme Rony, Adrian Galdran, Jose Dolz, and Ismail Ben Ayed. Class adaptive
713 network calibration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
714 Recognition (CVPR)*, pp. 16070–16079, June 2023a.
- 715
- 716 Bo Liu, Xingchao Liu, Xiaojie Jin, Peter Stone, and Qiang Liu. Conflict-averse gra-
717 dient descent for multi-task learning. In M. Ranzato, A. Beygelzimer, Y. Dauphin,
718 P.S. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Pro-
719 cessing Systems*, volume 34, pp. 18878–18890. Curran Associates, Inc., 2021. URL
720 [https://proceedings.neurips.cc/paper_files/paper/2021/file/
9d27fdf2477ffbff837d73ef7ae23db9-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/9d27fdf2477ffbff837d73ef7ae23db9-Paper.pdf).
- 721
- 722 Jiawei Liu, Changkun Ye, Ruikai Cui, and Nick Barnes. Self-calibrating vicinal risk minimisation
723 for model calibration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and
724 Pattern Recognition (CVPR)*, pp. 3335–3345, June 2024.
- 725
- 726 Shikun Liu, Edward Johns, and Andrew J. Davison. End-to-end multi-task learning with attention. In
727 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,
728 June 2019a.
- 729
- 730 Yang Liu, Shen Yan, Laura Leal-Taixé, James Hays, and Deva Ramanan. Soft augmentation for im-
731 age classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
732 Recognition (CVPR)*, pp. 16241–16250, June 2023b.
- 733
- 734 Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike
735 Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining
736 approach, 2019b.
- 737
- 738 Ze Liu, Han Hu, Yutong Lin, Zhuliang Yao, Zhenda Xie, Yixuan Wei, Jia Ning, Yue Cao, Zheng
739 Zhang, Li Dong, Furu Wei, and Baining Guo. Swin transformer v2: Scaling up capacity and
740 resolution. In *International Conference on Computer Vision and Pattern Recognition (CVPR)*,
741 2022b.
- 742
- 743 Jooyoung Moon, Jihyo Kim, Younghak Shin, and Sangheum Hwang. Confidence-aware learning
744 for deep neural networks. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th
745 International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning
746 Research*, pp. 7034–7044. PMLR, 13–18 Jul 2020. URL [https://proceedings.mlr.
747 press/v119/moon20a.html](https://proceedings.mlr.press/v119/moon20a.html).
- 748
- 749 Jishnu Mukhoti, Viveka Kulharia, Amartya Sanyal, Stuart Golodetz, Philip Torr, and Puneet
750 Dokania. Calibrating deep neural networks using focal loss. In H. Larochelle,
751 M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural In-
752 formation Processing Systems*, volume 33, pp. 15288–15299. Curran Associates, Inc.,
753 2020. URL [https://proceedings.neurips.cc/paper_files/paper/2020/
754 file/aeb7b30ef1d024a76f21a1d40e30c302-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/aeb7b30ef1d024a76f21a1d40e30c302-Paper.pdf).
- 755
- 752 Rafael Müller, Simon Kornblith, and Geoffrey E Hinton. When does label smoothing help? In
753 H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.),
754 *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.,
755 2019. URL [https://proceedings.neurips.cc/paper_files/paper/2019/
file/fl1748d6b0fd9d439f71450117eba2725-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/fl1748d6b0fd9d439f71450117eba2725-Paper.pdf).

- 756 Muhammad Akhtar Munir, Muhammad Haris Khan, Salman Khan, and Fahad Khan. Bridging
757 precision and confidence: A train-time loss for calibrating object detection. *IEEE Conference on*
758 *Computer Vision and Pattern Recognition (CVPR)*, 2023a.
- 759
760 Muhammad Akhtar Munir, Salman H Khan, Muhammad Haris Khan, Mohsen Ali, and Fa-
761 had Shahbaz Khan. Cal-detr: Calibrated detection transformer. In A. Oh, T. Nau-
762 mann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural*
763 *Information Processing Systems*, volume 36, pp. 71619–71631. Curran Associates, Inc.,
764 2023b. URL [https://proceedings.neurips.cc/paper_files/paper/2023/](https://proceedings.neurips.cc/paper_files/paper/2023/file/e271e30de7a2e462calf85cefa816380-Paper-Conference.pdf)
765 [file/e271e30de7a2e462calf85cefa816380-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/e271e30de7a2e462calf85cefa816380-Paper-Conference.pdf).
- 766 Mahdi Pakdaman Naeini, Gregory F. Cooper, and Milos Hauskrecht. Obtaining well calibrated
767 probabilities using bayesian binning. In *Proceedings of the Twenty-Ninth AAAI Conference on*
768 *Artificial Intelligence*, AAAI’15, pp. 2901–2907. AAAI Press, 2015. ISBN 0262511290.
- 769
770 Dexter Neo, Stefan Winkler, and Tsuhan Chen. Maxent loss: Constrained maximum entropy
771 for calibration under out-of-distribution shift. *Proceedings of the AAAI Conference on Arti-*
772 *ficial Intelligence*, 38(19):21463–21472, Mar. 2024. doi: 10.1609/aaai.v38i19.30143. URL
773 <https://ojs.aaai.org/index.php/AAAI/article/view/30143>.
- 774 Khanh Nguyen and Brendan O’Connor. Posterior calibration and exploratory analysis for natural
775 language processing models. In Lluís Màrquez, Chris Callison-Burch, and Jian Su (eds.), *Pro-*
776 *ceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pp.
777 1587–1598, Lisbon, Portugal, September 2015. Association for Computational Linguistics. doi:
778 10.18653/v1/D15-1182. URL <https://aclanthology.org/D15-1182>.
- 779 Jianmo Ni, Jiacheng Li, and Julian McAuley. Justifying recommendations using distantly-labeled
780 reviews and fine-grained aspects. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan (eds.),
781 *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and*
782 *the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp.
783 188–197, Hong Kong, China, November 2019. Association for Computational Linguistics. doi:
784 10.18653/v1/D19-1018. URL <https://aclanthology.org/D19-1018>.
- 785
786 Jeremy Nixon, Michael W. Dusenberry, Linchuan Zhang, Ghassen Jerfel, and Dustin Tran. Mea-
787 suring calibration in deep learning. In *Proceedings of the IEEE/CVF Conference on Computer*
788 *Vision and Pattern Recognition (CVPR) Workshops*, June 2019.
- 789 Jongyoun Noh, Hyekang Park, Junghyup Lee, and Bumsub Ham. Rankmixup: Ranking-based
790 mixup training for network calibration. In *Proceedings of the IEEE/CVF International Conference*
791 *on Computer Vision (ICCV)*, pp. 1358–1368, October 2023.
- 792 Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, D. Sculley, Sebastian Nowozin,
793 Joshua Dillon, Balaji Lakshminarayanan, and Jasper Snoek. Can you trust your model’s
794 uncertainty? evaluating predictive uncertainty under dataset shift. In H. Wallach,
795 H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett (eds.), *Ad-*
796 *vances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.,
797 2019. URL [https://proceedings.neurips.cc/paper_files/paper/2019/](https://proceedings.neurips.cc/paper_files/paper/2019/file/8558cb408c1d76621371888657d2eb1d-Paper.pdf)
798 [file/8558cb408c1d76621371888657d2eb1d-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/8558cb408c1d76621371888657d2eb1d-Paper.pdf).
- 799 Hyekang Park, Jongyoun Noh, Youngmin Oh, Donghyeon Baek, and Bumsub Ham. Acls: Adap-
800 tive and conditional label smoothing for network calibration. In *Proceedings of the IEEE/CVF*
801 *International Conference on Computer Vision (ICCV)*, pp. 3936–3945, October 2023.
- 802
803 Gabriel Pereyra, George Tucker, Jan Chorowski, Łukasz Kaiser, and Geoffrey Hinton. Regularizing
804 neural networks by penalizing confident output distributions. *arXiv*, abs/1701.06548, 2017.
- 805
806 John Platt and Nikos Karampatziakis. Probabilistic outputs for svms and comparisons to regularized
807 likelihood methods. In *Advances in Large Margin Classifiers*, 2007. URL [https://api.](https://api.semanticscholar.org/CorpusID:59806397)
808 [semanticscholar.org/CorpusID:59806397](https://api.semanticscholar.org/CorpusID:59806397).
- 809 Geoff Pleiss, Manish Raghavan, Felix Wu, Jon Kleinberg, and Kilian Q Weinberger. On fairness and
calibration. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan,

- 810 and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/b8b9c74ac526fffb2d39ab038d1cd7-Paper.pdf.
- 811
- 812
- 813
- 814 Rebecca Roelofs, Nicholas Cain, Jonathon Shlens, and Michael C. Mozer. Mitigating bias in calibration error estimation. In Gustau Camps-Valls, Francisco J. R. Ruiz, and Isabel Valera (eds.), *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume 151 of *Proceedings of Machine Learning Research*, pp. 4036–4054. PMLR, 28–30 Mar 2022. URL <https://proceedings.mlr.press/v151/roelofs22a.html>.
- 815
- 816
- 817
- 818
- 819 Hitesh Sapkota, Dingrong Wang, Zhiqiang Tao, and Qi Yu. Distributionally robust ensemble of lottery tickets towards calibrated sparse network training. In A. Oh, T. Neumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 62657–62681. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/c5cf13bfd3762821ef7607e63ee90075-Paper-Conference.pdf.
- 820
- 821
- 822
- 823
- 824
- 825 S. Schäffler, R. Schultz, and Klaus Weinzierl. Stochastic method for the solution of unconstrained vector optimization problems. *Journal of Optimization Theory and Applications*, 114:209–222, 01 2002. doi: 10.1023/A:1015472306888.
- 826
- 827
- 828
- 829 Ozan Sener and Vladlen Koltun. Multi-task learning as multi-objective optimization. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018. URL https://proceedings.neurips.cc/paper_files/paper/2018/file/432aca3ale345e339f35a30c8f65edce-Paper.pdf.
- 830
- 831
- 832
- 833
- 834 Zheyang Shen, Jiashuo Liu, Yue He, Xingxuan Zhang, Renzhe Xu, Han Yu, and Peng Cui. Towards out-of-distribution generalization: A survey. *ArXiv*, abs/2108.13624, 2021. URL <https://api.semanticscholar.org/CorpusID:237364121>.
- 835
- 836
- 837
- 838 Linwei Tao, Minjing Dong, and Chang Xu. Dual focal loss for calibration. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 33833–33849. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/tao23a.html>.
- 839
- 840
- 841
- 842
- 843 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 H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/file/36ad8b5f42db492827016448975cc22d-Paper.pdf.
- 844
- 845
- 846
- 847
- 848
- 849 Christian Tomani, Sebastian Gruber, Muhammed Ebrar Erdem, Daniel Cremers, and Florian Buettner. Post-hoc uncertainty calibration for domain drift scenarios. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 10124–10132, June 2021.
- 850
- 851
- 852
- 853
- 854
- 855
- 856
- 857
- 858
- 859
- 860
- 861
- 862
- 863
- Christian Tomani, Daniel Cremers, and Florian Buettner. Parameterized temperature scaling for boosting the expressive power in post-hoc uncertainty calibration. In Shai Avidan, Gabriel Brostow, Moustapha Cissé, Giovanni Maria Farinella, and Tal Hassner (eds.), *Computer Vision – ECCV 2022*, pp. 555–569, Cham, 2022. Springer Nature Switzerland. ISBN 978-3-031-19778-9.
- Christian Tomani, Futa Kai Waseda, Yuesong Shen, and Daniel Cremers. Beyond in-domain scenarios: Robust density-aware calibration. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 34344–34368. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/tomani23a.html>.
- Yoav Wald, Amir Feder, Daniel Greenfeld, and Uri Shalit. On calibration and out-of-domain generalization. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan

- 864 (eds.), *Advances in Neural Information Processing Systems*, volume 34, pp. 2215–2227. Curran Associates, Inc., 2021. URL https://proceedings.neurips.cc/paper_files/paper/2021/file/118bd558033a1016fcc82560c65cca5f-Paper.pdf.
- 865
866
867
- 868 Deng-Bao Wang, Lei Feng, and Min-Ling Zhang. Rethinking calibration of deep neural networks: Do not be afraid of overconfidence. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, volume 34, pp. 11809–11820. Curran Associates, Inc., 2021. URL https://proceedings.neurips.cc/paper_files/paper/2021/file/61f3a6dbc9120ea78ef75544826c814e-Paper.pdf.
- 869
870
871
872
873
- 874 Deng-Bao Wang, Lanqing Li, Peilin Zhao, Pheng-Ann Heng, and Min-Ling Zhang. On the pitfall of mixup for uncertainty calibration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 7609–7618, June 2023.
- 875
876
- 877 Jonathan Wenger, Hedvig Kjellström, and Rudolph Triebel. Non-parametric calibration for classification. In Silvia Chiappa and Roberto Calandra (eds.), *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pp. 178–190. PMLR, 26–28 Aug 2020. URL <https://proceedings.mlr.press/v108/wenger20a.html>.
- 878
879
880
881
- 882 Peiyao Xiao, Hao Ban, and Kaiyi Ji. Direction-oriented multi-objective learning: Simple and provable stochastic algorithms. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 4509–4533. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/0e5b96f97c1813bb75f6c28532c2ecc7-Paper-Conference.pdf.
- 883
884
885
886
887
- 888 Miao Xiong, Ailin Deng, Pang Wei W Koh, Jiaying Wu, Shen Li, Jianqing Xu, and Bryan Hooi. Proximity-informed calibration for deep neural networks. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 68511–68538. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/d826f5aadb26db488b8686097ceea2d1-Paper-Conference.pdf.
- 889
890
891
892
893
- 894 Jia-Qi Yang, De-Chuan Zhan, and Le Gan. Beyond probability partitions: Calibrating neural networks with semantic aware grouping. In A. Oh, T. Neumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 58448–58460. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/b693a240cf1009bff9fa4422141c9392-Paper-Conference.pdf.
- 895
896
897
898
899
- 900 Yaodong Yu, Stephen Bates, Yi Ma, and Michael Jordan. Robust calibration with multi-domain temperature scaling. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 27510–27523. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/b054fadflccd80b37d465f6082629934-Paper-Conference.pdf.
- 901
902
903
904
905
- 906 Amir R. Zamir, Alexander Sax, William Shen, Leonidas J. Guibas, Jitendra Malik, and Silvio Savarese. Taskonomy: Disentangling task transfer learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- 907
908
909
- 910 Hongyi Zhang, Moustapha Cissé, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. URL <https://openreview.net/forum?id=r1Ddpl-Rb>.
- 911
912
913
- 914 Jize Zhang, Bhavya Kailkhura, and T. Yong-Jin Han. Mix-n-match : Ensemble and compositional methods for uncertainty calibration in deep learning. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 11117–11128. PMLR, 13–18 Jul 2020. URL <https://proceedings.mlr.press/v119/zhang20k.html>.
- 915
916
917

918 Jun Zhang, Wen Yao, Xiaoqian Chen, and Ling Feng. Transferable post-hoc calibration on pre-
919 trained transformers in noisy text classification. *Proceedings of the AAAI Conference on Ar-*
920 *tificial Intelligence*, 37(11):13940–13948, Jun. 2023. doi: 10.1609/aaai.v37i11.26632. URL
921 <https://ojs.aaai.org/index.php/AAAI/article/view/26632>.
922

923 Linjun Zhang, Zhun Deng, Kenji Kawaguchi, and James Zou. When and how mixup improves
924 calibration. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu,
925 and Sivan Sabato (eds.), *Proceedings of the 39th International Conference on Machine Learning*,
926 volume 162 of *Proceedings of Machine Learning Research*, pp. 26135–26160. PMLR, 17–23 Jul
927 2022. URL <https://proceedings.mlr.press/v162/zhang22f.html>.

928 Yu Zhang and Qiang Yang. A survey on multi-task learning. *IEEE Transactions on Knowledge and*
929 *Data Engineering*, 34:5586–5609, 2017. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:11311635)
930 [CorpusID:11311635](https://api.semanticscholar.org/CorpusID:11311635).

931 Z.H. Zhou. *Ensemble Methods: Foundations and Algorithms*. CHAPMAN & HALL/CRC MA-
932 CHINE LEA. Taylor & Francis, 2012. ISBN 9781439830031. URL [https://books.](https://books.google.com.sg/books?id=MgR-wwEACAAJ)
933 [google.com.sg/books?id=MgR-wwEACAAJ](https://books.google.com.sg/books?id=MgR-wwEACAAJ).
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

A CALIBRATION ALGORITHMS AND METRICS

In this section, we further discuss in detail the various families of approaches commonly used to improve and measure neural network calibration.

A.1 ENTROPY-BASED METHODS

Entropy-based methods have played an important role in calibrating deep neural networks, as maximizing the entropy helps penalize overconfident predictions (Pereyra et al., 2017; Mukhoti et al., 2020; Neo et al., 2024). As mentioned in the main text, naively penalizing all predictions can cause underconfident predictions. While various works have proposed different approaches in controlling the entropy term, the Focal Loss (Lin et al., 2017; Mukhoti et al., 2020; Ghosh et al., 2022) and its variants offer adaptive/automated mechanisms in obtaining suitable values of γ for each sample.

While these automated mechanisms tend to help with ID calibration, many works fail to acknowledge the importance of OOD calibration since the parameters obtained during training/validation may not work during testing (Ovadia et al., 2019). As a work-around, we find that entropy-based methods can be extended to include OOD Maximum Entropy constraints (Jaynes, 1957; Neo et al., 2024) or Dual logit manipulation (Tao et al., 2023), showcasing the versatility of entropy-based methods. Since these methods all share the form of the Focal loss, we can easily pair all of them together into a single step.

A.2 REGULARIZERS

Mixup is an effective regularization technique that augments (Zhang et al., 2018) both input features and labels. Mixup works particularly well on both wider and deeper networks (Zhang et al., 2022) and can be particularly useful in improving network calibration (Thulasidasan et al., 2019; Chidambaram & Ge, 2024). As an extension to vanilla Mixup, RankMixup (Noh et al., 2023) can be used to ensure that the augmented samples have lower confidences than the original samples.

A.3 MARGIN-BASED METHODS

Margin-based methods tend to restrict model confidences by a constant margin/factor. For example, label smoothing (LS) (Müller et al., 2019) softens the targets using a constant factor ϵ . Mathematically, the smoothed label s_i is acquired after uniformly adjusting the target $s_i = (1 - \epsilon)y_k + \frac{\epsilon}{K}$, which is then used to train the network. Although vanilla LS can be used to improve miscalibration, imposing a constant smoothing factor for all training labels can lead to under-confident predictions. Furthermore, searching for a suitable ϵ is computationally expensive as it requires a grid-search across multiple models during the training phase.

Instead of implementing a fixed constant, several works have been proposed to adaptively or conditionally approximate the label smoothing function during training. For example, MDCA (Hebbalaguppe et al., 2022) utilizes a regularization term, which enforces predicted confidences to be as close to the average accuracy as possible. This can lead to a parabolic smoothing function (Park et al., 2023), that is adaptively dependent on the predicted confidences. Which can be problematic, since both high and low confidence predictions are weakly penalized. Another approach would be to only conditionally smooth predictions based on a margin. For instance, MBLS (Liu et al., 2022a) and CALS-ALM (Liu et al., 2023a) propose to restrict output logits by a user defined margin, but can be sensitive to hyper parameter settings. CRL (Moon et al., 2020) ordinarily ranks predictions based on the number of times each sample is predicted correctly, however it requires a buffer to store the correctness history. Which can be empty during the earlier stages of training and idle during later phases when the model’s accuracy is high.

By adopting a smoothing and indicator function, the Adaptive Conditional Label Smoothing (ACLS) (Park et al., 2023) method seeks to combine the benefits of both adaptive and conditional methods without the use of an additional correctness history.

1026 A.4 POST-HOC PROCESSING

1027
1028 The fundamental idea behind post-hoc processing methods is to obtain a mapping func-
1029 tion/temperature that modifies the model’s logits thus changing it’s predicted confidence. The most
1030 popular post-processing step is the vanilla temperature scaling (TS) (Platt & Karampatziakis, 2007),
1031 which manipulates the model’s confidences without changing the final class label predictions. For
1032 example, a value of $T < 1$ leads to a lower entropy or “peaky” distributions and a value of $T > 1$
1033 gives higher entropy or “flatter” predictions.

1034 The typical approach in obtaining the temperature parameter, is to minimize the *average* calibration
1035 error or NLL over a separate validation set. While vanilla TS has been found to be effective in
1036 reducing network over-confidence (Guo et al., 2017), it generally reduces the confidence of every
1037 sample - even when predictions are correct. Other forms of post-hoc processing include calibration
1038 using, model ensembles (Zhang et al., 2020), splines (Gupta et al., 2021) and distribution matching
1039 (Kuleshov & Deshpande, 2022; Tomani et al., 2023). For our post-processing step, we use AdaTS
1040 since it is the SOTA method for post-processing methods and adaptively chooses a samplewise
1041 temperature for scaling model predictions.

1042 A.5 CALIBRATION METRICS

1043
1044 **Expected Calibration Error (ECE):** The ECE is the most widely used metric in the literature
1045 and directly tied to the definition of calibration (Guo et al., 2017; Tomani et al., 2021). By splitting
1046 the predicted confidences in B evenly separated bins, each containing n_b samples. The ECE is then
1047 simply a scalar measuring the weighted errors between the *acc* and *conf* of each bin (Naeini et al.,
1048 2015): $ECE = \sum_{b=1}^B \frac{n_b}{N} |acc(b) - conf(b)|$. Despite the ECE’s popularity, many recent works have
1049 pointed out the limitations of the ECE, such as bin size sensitivity and it’s lack of consideration for
1050 classwise calibration. For a fair and thorough analysis, we introduce other calibration metrics that
1051 cover the weaknesses of the ECE.
1052

1053
1054 **Classwise ECE (CECE):** As most calibration metrics typically only considers the max confi-
1055 dence probabilities, the CECE considers the macro-averaged ECE of all K classes. Predictions are
1056 binned individually for each respective class and the calibration error is measured for each class
1057 level bin (Nixon et al., 2019). $CECE = \frac{1}{K} \sum_{b=1}^B \sum_{k=1}^K \frac{n_{b,k}}{N} |acc(b, k) - conf(b, k)|$.

1058
1059 **Overconfidence Error (OE):** For safety-critical applications, overconfident mispredictions are
1060 potentially hazardous. The OE penalizes overconfident bins that have higher confidences than accu-
1061 racy (Thulasidasan et al., 2019): $OE = \sum_{b=1}^B \frac{n_b}{N} [conf(b) \times \max(conf(b) - acc(b), 0)]$.

1062
1063 **Kolmogorov-Smirnov Error (KSE):** As many calibration metrics are often sensitive to the num-
1064 ber of B bins used during the partitioning of empirical distributions. The KSE (Gupta et al., 2021)
1065 is a bin-free alternative that numerically approximates the differences between two empirical cumu-
1066 lative distributions. The KSE for top-1 classification is given as the following integral, with z_k
1067 denoting the predicted probabilities: $KSE = \int_0^1 |P(k|z_k) - z_k| P(z_k) dz_k$.

1068
1069 **Adaptive ECE (AdaECE):** as the ECE is known to be biased towards higher confidence bins,
1070 the AdaECE (Nguyen & O’Connor, 2015) is proposed to adaptively/evenly measure samples across
1071 bins: $AdaECE = \sum_{b=1}^B \frac{n_b}{N} |acc(b) - conf(b)|$ s.t. $\forall b, i \cdot |B_b| = |B_i|$.

1072
1073 **Negative Log-likelihood (NLL):** Commonly referred to as cross entropy in deep learning. The
1074 NLL (Hastie et al., 2001) measures the alignment between a model’s confidence $P_i(y_k|x)$ and targets
1075 y_k : $NLL = -\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K y_k \log P_i(y_k|x)$.

Algorithm 1: Peacock - Unified Multi-Objective Optimization Calibration Framework

```

1080 Data: Given training and validation set  $D_{\text{train}} = (x_i, y_i)_{i=1}^N$ ,  $D_{\text{val}} = (x_v, y_v)_{v=1}^V$ 
1081
1082 1: Initialize neural network parameters  $\theta$ , learning rate schedule  $\eta$  and uniformly distributed weights  $w_t = \frac{1}{A}$ 
1083 2: Compute the global and local expectations for the mean and variance constraints  $\mu, \sigma^2$ 
1084 3:  $\hookrightarrow \mathbb{E}[\mathcal{Y}] = \mu$  and  $\mathbb{E}[\mathcal{Y}^2] = \sigma^2$ 
1085 4: Solve numerically for  $\lambda_\mu \leftarrow \text{NewtonRaphson}()$  // MaxEnt loss root-finder
1086 5: for  $e \in \text{epochs}$  do
1087 6:   for  $i \in B$  do // Sample mini-batch of size  $B$ 
1088 7:     Perform FastMixup on images:  $\tilde{x} = \beta x_i + (1 - \beta)x_j$  // RankMixup
1089 8:     Perform FastMixup on labels:  $\tilde{y} = \beta y_i + (1 - \beta)y_j$  // RankMixup
1090 9:     Compute 1v1 loss:  $\mathcal{L}_{\text{CPC}}^{\text{1v1}} = -\frac{1}{(C-1)} \sum_{j \neq y} \log \frac{P_i y}{P_i y + P_i j}$  // CPC loss
1091 10:    if  $\gamma_{t,b} \geq 0$  then
1092 11:       $\mathcal{L}_{\text{AdaFocal}}^{\text{Dual}} = -\sum_k (1 - P_i + P_j)^{\gamma_{t,b}} \log P_i$  // Dual AdaFocal loss
1093 12:    else if  $\gamma_{t,b} < 0$  then
1094 13:       $\mathcal{L}_{\text{AdaFocal}}^{\text{Dual}} = -\sum_k (1 + P_i + P_j)^{|\gamma_{t,b}|} \log P_i$  // Inverse Dual AdaFocal loss
1095 14:    Compute MaxE loss  $\mathcal{L}_{\text{ME}} = \lambda_\mu (\sum_k f(\mathcal{Y}) P_i(y_k|x) - \mu_G + \sum_k f(\mathcal{Y}) P_i(y_k|x) - \mu_{Lk})$  // MaxEnt loss
1096 15:    Compute MRL loss  $\mathcal{L}_{\text{MRL}} = \max(0, \max_k \hat{P} - \max_k P + m_{\text{MRL}})$  // RankMixup
1097 16:    if  $j = \tilde{y}$  then
1098 17:       $\mathcal{L}_{\text{ACLS}} = \lambda_1 \max(0, g_j^\theta(x) - \min_k (g_k^\theta(x)) - m_{\text{ACLS}})^2$  // ACLS regularizer
1099 18:    else if  $j \neq \tilde{y}$  then
1100 19:       $\mathcal{L}_{\text{ACLS}} = \lambda_2 \max(0, g_{\tilde{y}}^\theta(x) - g_j^\theta(x) - m_{\text{ACLS}})^2$  // ACLS regularizer
1101 20:     $w_t = \text{ImportancePeacock}(\mathcal{L}_t(\theta))$  // Compute importance loss weights
1102 21:
1103 22:    Compute Peacock:
1104 23:     $\hookrightarrow \mathcal{L}_{\text{Peacock}} = \mathcal{L}_{\text{AdaFocal}}^{\text{Dual}} + w_1 \mathcal{L}_{\text{constraints}}^{\text{ME}} + w_2 \mathcal{L}_{\text{CPC}}^{\text{1v1}} + w_3 \mathcal{L}_{\text{MRL}}$ 
1105 24:     $\theta_{\text{new}} \leftarrow \theta_{\text{old}} - \eta \nabla_{\theta} \mathcal{L}_{\text{Peacock}}$  // Update parameters  $\theta$  by gradient descent
1106 25:  return  $\theta$ 
1107 26:
1108 27: Apply temperature scaling:  $\theta_{\text{AdaTS}} \leftarrow \text{AdaptiveTS}(D_{\text{val}}, \theta)$  // AdaTS
1109 28: Function  $\text{NewtonRaphson}()$ :
1110 29:    $\delta = 1e-15$  // A small tolerance or stopping condition
1111 30:   while  $g(\lambda) > \delta$  do
1112 31:      $\lambda_{n+1} = \lambda_n - \frac{g(\lambda)}{g'(\lambda)}$  // Update Lagrange Multipliers  $\lambda_n$ 
1113 32:   return  $\lambda_n$ 
1114 33: Function  $\text{ImportancePeacock}()$ :
1115 34:    $\min_{w_t} \left\{ \left\| \sum_{t=1}^A w_t \sqrt{\frac{\Delta_{\theta} \mathcal{L}_t(\theta)}{\eta}} \right\|_2^2 \mid \sum_{t=1}^A w_t = 1, w_t \geq 0 \quad \forall t \right\}$ 
1116 35:   return  $w_t$ 
1117 36: Function  $\text{AdaptiveTS}(D_{\text{val}}, \theta)$ :
1118 37:   Initialize VAE and MLP parameters  $Q, \phi$ 
1119 38:   while  $t < \text{steps}$  do
1120 39:     for  $v \in B$  do // Sample mini-batch of size  $B$ 
1121 40:        $\nabla_{VAE} \leftarrow \nabla_{\text{ELBO}}[\Phi(x)]$ 
1122 41:        $\tilde{q} = \{\log P(z|y) | \forall y\}$   $z \sim Q_{\phi}(z|x)$ 
1123 42:        $\nabla_T \leftarrow \log(\text{softmax}(g_{\theta}/T))$ 
1124 43:        $(\theta, \phi)_{t+1} \leftarrow (\theta, \phi)_t - \alpha_{\text{tr}}(\nabla_{VAE} + \nabla_T)$ 
1125 44:   return  $\theta_{\text{AdaTS}}$ 

```

B IMPLEMENTATION DETAILS FOR PEACOCK**B.1 ALGORITHM DETAILS AND HYPERPARAMETERS**

For our implementation of *Peacock*, we first select the mean constraint for of MaxEnt loss as our starting algorithm and compute Lagrange multipliers λ_n using the Newton Raphson method. This step is performed in $\mathcal{O}(n)$ time using the helper function $g(\lambda)$ and its derivative $g'(\lambda)$ before model training begins.

For each iteration, the pairwise 1v1 constraints of CPC loss are first computed before incorporating the adaptive γ selection mechanism of AdaFocal loss. This step also includes the second highest confidence $P_i(y_j|x)$ from Dual Focal loss to AdaFocal loss. This way we can reduce compute overhead by combining both calibration methods into a single step: $\mathcal{L}_{\text{AdaFocal}}^{\text{Dual}} = \mathcal{L}_{\text{AdaFocal}} + \mathcal{L}_{\text{Dual}}$.

Next, RankMixup is performed for every sampled input image and label, with the MRL loss computed using the coefficients α and m . For texts datasets, we perform RankMixup at the feature level. Although vanilla Mixup has been found to hurt ID calibration performance (Wang et al., 2023), our findings suggest that by combining RankMixup with other algorithms good balance between ID and OOD calibration can be achieved. In our experience, we find that a large ACLS margin, can lead

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145

Hyperparameters	Values
Learning rate η	0.1
Batch size	512 or 256
Optimizer	SGD or Adam
Scheduler	Cosine Annealing or Fixed
Epochs	200 or 50
Margin m_{ACLS}	6.0
Mixup α	1.0
Mixup margin m_{MRL}	2.0
γ starting	1.0
γ max	20.0
γ min	-2.0
No. of bins B	15.0
Learning rate for attention block	3e-4

1146
1147
1148

Table 4: Hyperparameters used for optimizing *Peacock*

1149
1150
1151
1152
1153
1154
1155
1156

to numerical instability when the number of classes is large, thus we fixed $m_{\text{ACLS}} = 6.0$. For completeness, we include the ACLS step in Algorithm 1, however in our ablation study we show that ACLS does not improve overall calibration performance and is not included during optimization or our final proposed version of *Peacock*. Next, an optional post-processing step using AdaTS is performed by learning the adaptive temperature on a separate validation set. Finally, for the importance weighted form of *Peacock*, we randomly initialize a self-attention block and optimize it with Eq. (16) with Adam optimizer and a learning rate of $3e-4$ to learn a set of importance weights for each loss term.

1157
1158
1159
1160
1161
1162
1163

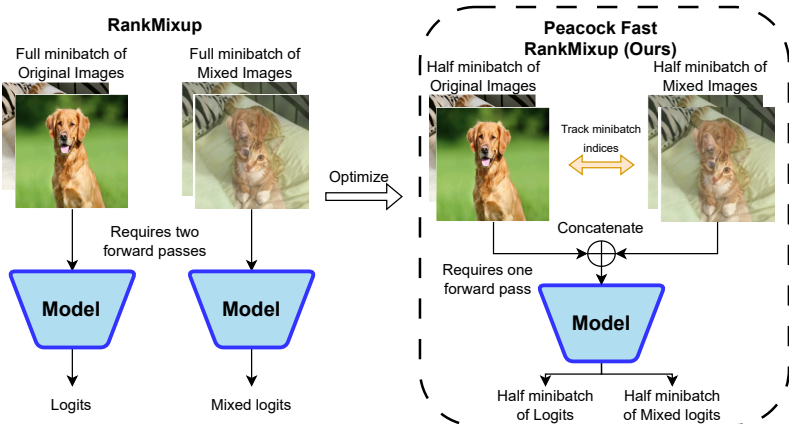
Hyperparameters In general, we try to keep the default settings of each algorithm. However, when trying to combine multiple of these components, it may become inevitable for some tuning to be performed. Indeed, performing a grid-search would be the best way to obtain the optimal hyperparameters. However, as discussed in our *Limitations*, the number of parameters scale exponentially with the number of calibration components selected for optimization. This can be easily become very compute intensive and would not be the focus of our work.

1164 B.2 ACCELERATING RANKMIXUP

1166
1167
1168
1169

Fig. 6 illustrates the comparisons between the original RankMixup method and the optimized version proposed in our paper. RankMixup, in its original form, requires two forward passes during training: one for a full minibatch (e.g., 512) of original images and another full minibatch of mixed images. This process can be computationally expensive, especially for large datasets or complex

1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184



1185
1186
1187

Figure 6: During training, RankMixup requires two forward passes: one for original images and one for mixed images, in order to compute \mathcal{L}_{MRL} . We optimize RankMixup by mixing images and labels batchwise, resulting in a 2x speed up during training.

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241

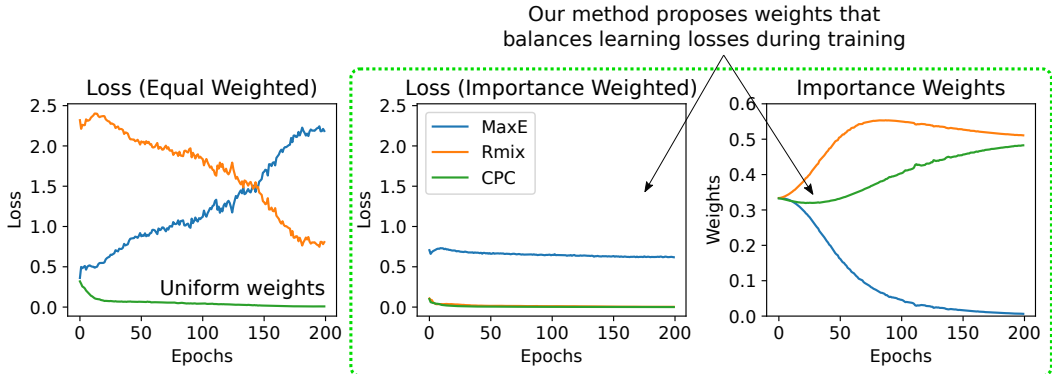


Figure 7: When using equal weights (left), the model optimizes loss terms equally. Our direction-weighted self-attention block, on the other hand, learns to dynamically adjust the importance of each loss term during training, enabling a more balanced optimization of the overall objective. All weights are initialized uniformly with plots smoothed for readability.

models. As a workaround, we propose an optimized variant of FastRankMixup, which addresses this limitation by dividing a full batch of images into two halves: containing a minibatch of half original and half mixed images (e.g., $512 \div 2 = 256$).

This way, we only require a single forward pass instead of the two forward passes, delivering a 2x speedup during training compared to the original RankMixup implementation. This improvement in training efficiency can be particularly beneficial for large-scale training tasks, where computational resources are often constrained. A caveat to this method is that the minimum batchsize required will always be two, as at least two samples are needed.

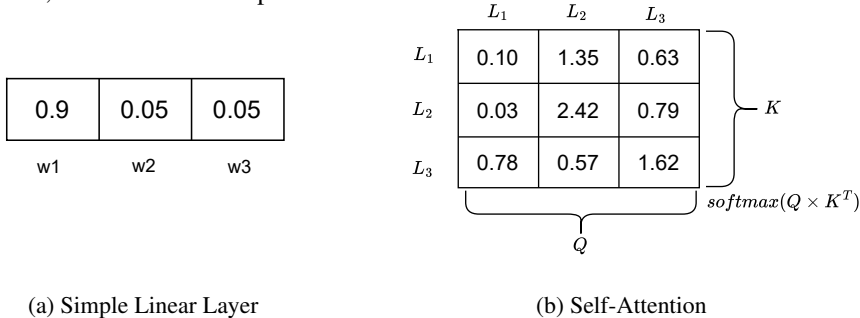


Figure 8: Comparisons between a simple linear layer versus self-attention for *Peacock*. Self-attention is better suited for capturing the complex relationships between loss terms.

B.3 LEARNING DYNAMICS OF DIRECTION WEIGHTED SELF-ATTENTION BLOCK

The importance-weighted formulation of *Peacock* utilizes a novel direction weighted self-attention block. This subsection discusses the learning dynamics and certain key considerations of the self-attention block. Fig. 7 illustrates the differences between the learning dynamics of the equal and importance weighted *Peacock* on CIFAR100. By assigning equal weights to each loss term, the model regards all auxiliary losses equally. Conversely, our proposed direction weighted self-attention block outputs importance weights at every timestep, using Eq. (16). This leads to an overall balanced and more stable learning process during optimization. Note that all loss terms are normalized before being passed into the self-attention block during training. Additionally, the direction weighted self-attention block provides certain key benefits:

- **Softmax of Self-Attention:** The softmax function of the self-attention block implicitly enforces KKT conditions, simplifying the optimization process.
- **Better learns relationships across losses:** The design of self-attention enables better learning of inter-dependencies among loss terms, compared to a linear layer (see Fig. 8).

C PROOFS

C.1 TEMPERATURE-SCALED BOUNDS

Consider a temperature/mapping function T which scales the output logits/hypothesis h^θ of a model. Then the average of each temperature scaled hypothesis is given as:

$$\overline{\text{ECE}}^2(\mathcal{T}(h^\theta)) = \frac{1}{A} \sum_{t=1}^A \text{ECE}^2(\mathcal{T}(h_t^\theta)) = \frac{\text{ECE}(\mathcal{T}(h_1^\theta))^2 + \text{ECE}(\mathcal{T}(h_2^\theta))^2 + \dots + \text{ECE}(\mathcal{T}(h_t^\theta))^2}{A} \quad (17)$$

Considering equal contributions of each individual temperature scaled hypothesis, the temperature scaled multi-objective learner $\mathcal{T}(H^\theta)$ has the expected squared ECE:

$$\text{ECE}^2(\mathcal{T}(H^\theta)) = \mathbb{E} \left[\left(\frac{1}{A} \sum_{t=1}^A \text{ECE}(\mathcal{T}(h_t^\theta)) \right)^2 \right] = \int \left(\frac{1}{A} \sum_{t=1}^A \text{ECE}(\mathcal{T}(h_t^\theta)) \right)^2 p(x) dx \quad (18)$$

which follows the same bounds as previously defined in the main paper.

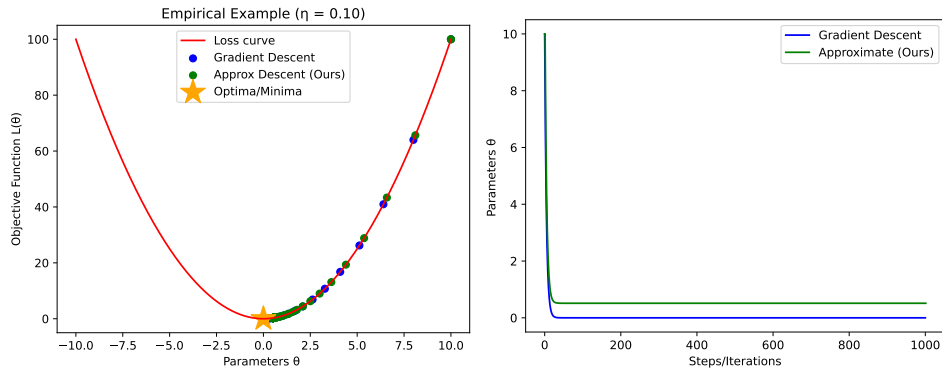
$$\text{ECE}^2(\mathcal{T}(H^\theta)) \leq \overline{\text{ECE}}^2(\mathcal{T}(h^\theta)) \quad (19)$$

Empirically, Table 2 demonstrates that if the same mapping function or temperature \mathcal{T} is applied to each hypothesis (e.g., AdaTS), then the average of the scaled combined learner will also obey the upper bound of the above inequality.

C.2 ESTIMATING THE GRADIENT

Recall in Section 4.2 of our main paper, the direct computation of $\nabla_{\theta} \mathcal{L}_t(\theta)$ requires the use of retaining the computational graph⁴ after the backward pass, which can be compute intensive and significantly slows down training time. In this section, we demonstrate that decrease rate estimates for each loss term can act as alternatives to direct gradient recomputation. By simply storing the previous loss value computed (single step look-back), we can avoid graph retention during the optimization for w_t . Using a simple example, we also show that our decrease rate estimates are closely related to the solutions obtained using gradient descent. For simplicity, we denote the partial derivatives as $\nabla_{\theta} \mathcal{L}_t(\theta) = \frac{\partial \mathcal{L}_t(\theta)}{\partial(\theta)}$ and the difference between old and new parameters as $\Delta_{\theta} \mathcal{L}_t(\theta)$.

⁴ For more details, see Pytorch autograd framework: <https://pytorch.org/docs/stable/autograd.html>



(a) Toy sketch illustrating the solution of our method (green) compared to gradient descent (blue). Our method closely approximates gradient descent.

(b) We compare the solutions given by our method and gradient descent, for $\eta = 0.1$, our solution is close to the solution by gradient descent.

Consider the following task of approximating $\nabla_{\theta}\mathcal{L}_t(\theta)$. By using the first order form of Taylor’s Theorem, the loss gradients can be rewritten as the following equation:

$$\nabla_{\theta}\mathcal{L}_t(\theta) = \frac{\mathcal{L}_t(\theta_{\text{new}}) - \mathcal{L}_t(\theta_{\text{old}})}{\Delta\theta} + \epsilon(\theta) = \frac{\Delta_{\theta}\mathcal{L}_t(\theta)}{\Delta\theta} + \epsilon(\theta) \quad (20)$$

where $\Delta_{\theta}\mathcal{L}_t(\theta)$ is the rate of change for each loss term with respect to the change of model parameters θ_{new} and θ_{old} , paired by a small error term $\epsilon(\theta)$. From the gradient descent update rule, the change in model parameters is given by:

$$\begin{aligned} \theta_{\text{new}} &= \theta_{\text{old}} - \eta\nabla_{\theta}\mathcal{L}_t(\theta) \\ \Delta\theta &= -\eta\nabla_{\theta}\mathcal{L}_t(\theta) \end{aligned} \quad (21)$$

where the difference between new and old network parameters are obtained using the gradients and a learning rate η . By substituting Eq. (21) into Eq. (20):

$$\begin{aligned} \nabla_{\theta}\mathcal{L}_t(\theta)^2 &= \frac{\Delta_{\theta}\mathcal{L}_t(\theta)}{-\eta} + \epsilon(\theta) = \frac{\mathcal{L}_t(\theta_{\text{old}}) - \mathcal{L}_t(\theta_{\text{new}})}{\eta} + \epsilon(\theta) \\ & \text{*Note the flip in sign} \\ \nabla_{\theta}\mathcal{L}_t(\theta) &= \sqrt{\frac{\Delta_{\theta}\mathcal{L}_t(\theta)}{\eta} + \epsilon(\theta)} \approx \sqrt{\frac{\Delta_{\theta}\mathcal{L}_t(\theta)}{\eta}} \end{aligned} \quad (22)$$

with the small error term $\epsilon(\theta)$ dropped.

Key Assumptions: Our main paper highlighted the essential assumptions underlying this formulation: 1.) The loss terms \mathcal{L}_t are convex and optimizable by gradient descent. 2.) Each loss term monotonically decreases i.e., the loss evaluated at previous iterations will always be strictly larger than the loss at the current iteration $\mathcal{L}_t(\theta_{\text{old}}) > \mathcal{L}_t(\theta_{\text{new}})$. This assumption ensures that the ratio $\sqrt{\frac{\Delta_{\theta}\mathcal{L}_t(\theta)}{\eta}}$ remains positive, avoiding the computation of complex numbers. Moreover, the small learning rates commonly used in deep learning frameworks tend to be sufficiently small (e.g., $\eta = 2.5\text{e-}4$) allowing for accurate linear approximations. In practice, we can apply the ReLU function to the gradient update, i.e $\sqrt{\text{ReLU}(\frac{\Delta_{\theta}\mathcal{L}_t(\theta)}{\eta})}$ if the gradient descent step leads to an increase in the loss, violating the assumption that $\mathcal{L}_t(\theta_{\text{old}}) > \mathcal{L}_t(\theta_{\text{new}})$. This ensures that the update is scaled down or ignored in the optimization process.

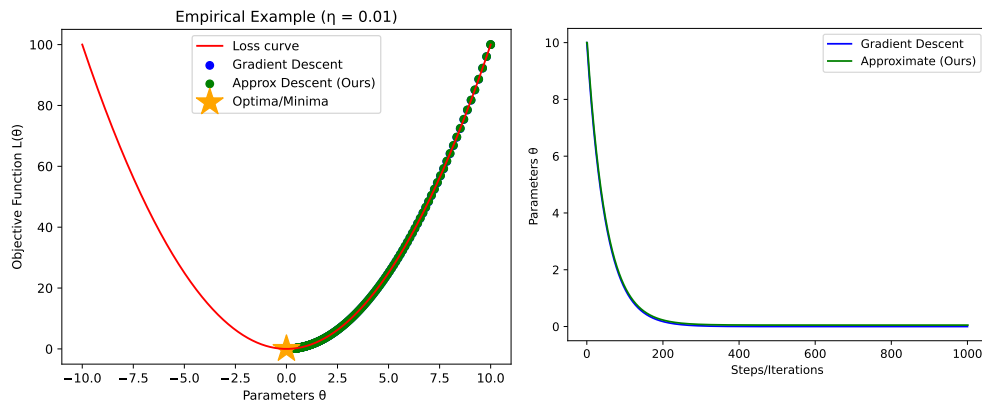
Simple Empirical Example We further support our findings by including a simple empirical example comparing gradient descent and our proposed method. Consider a smooth, convex objective function $\mathcal{L}_{\theta} = \theta^2$. Fig. 9a illustrates our goal of obtaining a set of parameters θ such that \mathcal{L}_{θ} is minimized. For a fixed learning rate of $\eta = 0.1$, 1000 iteration steps and a starting point of $\theta = 10$, our solution given by our method (in green) is relatively close compared to the solution given by gradient descent (in blue), with a slight delay and an error of roughly 0.5. We can further improve our method’s solution by reducing the learning rate to $\eta = 0.01$, which provides an even closer estimate to the solutions given by gradient descent and a reduced relative error of roughly 0.05.

D SUPPLEMENTARY EXPERIMENTS AND RESULTS

D.1 DATASET DETAILS

Synthetic OOD We train our models with clean images from the original CIFAR, TinyImageNet and evaluate their OOD performance on their corrupted forms CIFAR-C, TinyImageNet-C.

1. CIFAR10/CIFAR100 (Krizhevsky & Hinton, 2009) RGB images of size (32x32) containing ten and hundred classes. The training/validation/testing sets contain 45,000/5,000/10,000 samples respectively.



(a) By reducing the learning rate, our method (green) provides an even closer estimate to gradient descent (blue).

(b) We compare the solutions given by our method and gradient descent for $\eta = 0.01$, our method can be improved by reducing the learning rate.

2. TinyImagenet (Deng et al., 2009) A miniature version of the ImageNet dataset containing images of size (64x64) of 200 classes. There are 100,000 images for training and 10,000 images for validation/testing.
3. CIFAR10-C/CIFAR100-C/TinyImagenet-C (Hendrycks & Dietterich, 2019) A widely popular calibration benchmark, containing corrupted variants of CIFAR and TinyImageNet. Standard image corruptions (total of 19) are applied on the original test sets across five increasing levels of severities.

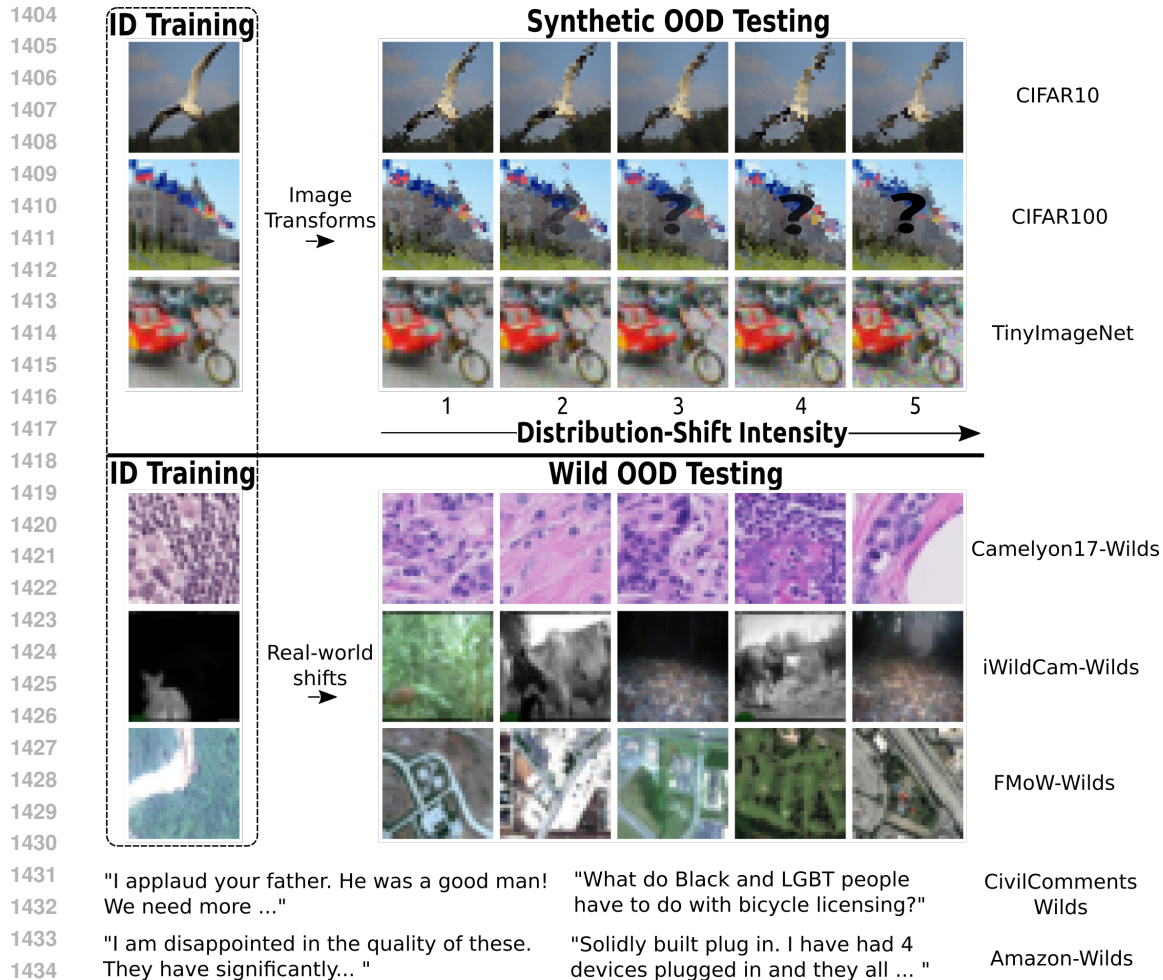
Real-world OOD For wild OOD, we learn our models using the provided ID training sets and OOD sets for validation and testing (Koh et al., 2021).

1. Camelyon17 (Bandi et al., 2019): A binary task to detect if a (32x32) cell tissue slide is benign or malignant. The images are collected across different hospitals with equipment that may vary OOD from the training set.
2. iWildCam (Beery et al., 2020): Animal species tend to vary across different backgrounds and terrains. The goal is to classify 182 animal classes collected from camera traps deployed in different areas of the wilderness.
3. FMoW (Christie et al., 2018): Satellite imagery of topographies and buildings alike tend to differ greatly across countries. The task is classify the OOD shifted terrains from one out of 62 classes.
4. CivilComments (Borkan et al., 2019): A binary text-classification task, where the model needs to identify toxic comments. The OOD shifts stem from inputs collected from differing demographics such as gender, religion, etc.
5. Amazon (Ni et al., 2019): A consumer-rating dataset where the input is a text review, with a label from a 1-to-5 star rating.

D.2 SUPPLEMENTARY EXPERIMENTS AND RESULTS

ID Results: As demonstrated in Table 5, our synthetic benchmark results confirm *Peacock*'s highly competitive performance on ID test sets. As discussed in the main text, *Peacock*'s calibration error is inherently bounded by the average calibration error of its constituent components. Consequently, even if some components underperform, *Peacock*'s overall calibration remains well-calibrated, irrespective of whether the data is in-distribution (ID) or out-of-distribution (OOD).

Additional OOD Results: Table 5, shows OOD supplementary results evaluated using AdaECE (Nguyen & O'Connor, 2015) and OE (Thulasidasan et al., 2019). Our analysis using these additional metrics aligns with the results presented in our primary findings. While the theoretical proofs in our



1438 Figure 11: Covariate shifts can be simulated using common image corruptions or caused by natural differences during data collection in-the-wild.

1439

1440

1441

1442

1443

1444

1445

1446

Dataset	Metric	MaxEnt	AdaFocal	RankMixup	CPC	Dual	ACLS	Peacock (Eq.)	Peacock (Impt.)
CIFAR10	Acc. \uparrow	94.0 \pm 0.2	93.4 \pm 0.4	94.1 \pm 0.1	94.6 \pm 0.1	94.5 \pm 0.1	94.3 \pm 0.1	93.8 \pm 0.1	93.9 \pm 0.2
	ECE \downarrow	1.1 \pm 0.1	0.8 \pm 0.1	3.1 \pm 0.1	2.9 \pm 0.1	1.3 \pm 0.1	3.0 \pm 0.1	0.6 \pm 0.1	0.6 \pm 0.1
	CECE \downarrow	0.4 \pm 0.1	0.3 \pm 0.1	2.8 \pm 0.1	2.7 \pm 0.1	0.4 \pm 0.1	2.7 \pm 0.1	0.2 \pm 0.1	0.2 \pm 0.1
	NLL \downarrow	249.8 \pm 0.4	232.7 \pm 0.1	346.9 \pm 0.2	394.4 \pm 0.2	253.7 \pm 0.1	345.3 \pm 0.3	224.4 \pm 0.4	224.5 \pm 4.1
CIFAR100	Acc. \uparrow	73.8 \pm 0.1	75.8 \pm 0.1	74.9 \pm 0.3	74.7 \pm 0.1	75.4 \pm 0.1	75.3 \pm 0.1	74.5 \pm 0.4	73.5 \pm 0.5
	ECE \downarrow	5.4 \pm 0.5	6.8 \pm 0.1	4.9 \pm 0.1	8.8 \pm 0.3	9.1 \pm 0.1	4.5 \pm 0.3	4.1 \pm 0.3	3.9 \pm 0.2
	CECE \downarrow	0.2 \pm 0.1	0.1 \pm 0.1	2.9 \pm 0.1	2.3 \pm 0.2	0.1 \pm 0.1	1.7 \pm 0.1	0.1 \pm 0.1	0.1 \pm 0.1
	NLL \downarrow	312.7 \pm 0.6	306.5 \pm 2.3	348.6 \pm 0.7	432.2 \pm 0.8	319.8 \pm 0.4	346.6 \pm 1.0	298.3 \pm 1.2	283.9 \pm 0.1
TinyImageNet	Acc. \uparrow	63.1 \pm 0.3	60.8 \pm 0.1	61.6 \pm 0.3	65.0 \pm 0.3	63.2 \pm 0.3	64.9 \pm 0.1	61.2 \pm 0.1	62.3 \pm 0.4
	ECE \downarrow	18.2 \pm 0.3	6.1 \pm 0.5	5.5 \pm 0.3	10.3 \pm 0.4	6.8 \pm 0.1	5.0 \pm 0.3	6.2 \pm 0.3	3.9 \pm 0.3
	CECE \downarrow	0.1 \pm 0.1	0.1 \pm 0.1	0.1 \pm 0.1	0.1 \pm 0.1	0.1 \pm 0.1	0.1 \pm 0.1	0.1 \pm 0.1	0.1 \pm 0.1
	NLL \downarrow	322.0 \pm 0.5	320.5 \pm 0.3	343.2 \pm 1.0	358.5 \pm 1.6	339.2 \pm 1.0	342.3 \pm 0.4	333.9 \pm 2	324.2 \pm 2.4

1448 Table 5: We report the ID test scores (%) for reruns computed across 3 seeds for *Peacock* and its

1449 components.

1450

1451

1452 main paper and Appendix C are explicitly stated only for the ECE, we anticipate that our arguments

1453 remain valid for other calibration metrics, which are often derivatives or closely related to ECE. We

1454 intend to explore this aspect further in future research.

1455

1456 **Additional Multi-Objective Optimization Results:** Additional results for our proposed

1457 weighted-importance formulation are provided in Table 8 and Table 9. Our results highlight the

versatility, effectiveness and speed across a wide variety of different architectures and methods.

Dataset	Metric	MaxEnt	AdaFocal	RankMixup	CPC	Dual	ACLS	Peacock (Eq.)	Peacock (Impt.)
CIFAR10-C	AdaECE ↓	6.9±0.1	6.2±0.4	11.5±0.2	10.7±0.4	7.2±0.2	11.5±0.1	6.3±0.1	6.2±0.3
	OE ↓	3.9±0.3	3.0±0.3	9.6±0.2	9.1±0.4	3.8±0.1	9.5±0.1	3.3±0.1	3.5±0.2
CIFAR100-C	AdaECE ↓	11.0±0.1	13.7±0.3	8.4±0.2	13.7±0.2	15.5±0.1	10.2±0.3	9.7±0.3	9.6±0.3
	OE ↓	0.5±0.1	0.7±0.1	2.5±0.2	1.8±0.1	0.7±0.1	2.01±0.1	1.3±0.1	1.6±0.1
TinyImageNet-C	AdaECE ↓	12.6±0.3	13.9±0.3	20.2±0.2	16.3±0.2	18.8±0.2	20.6±0.4	10.4±0.2	10.7±0.2
	OE ↓	4.4±0.3	4.5±0.3	10.4±0.2	8.5±0.2	9.4±0.2	11.1±0.4	2.3±0.2	2.3±0.2
Camelyon17	AdaECE ↓	12.3±0.4	20.4±0.1	22.4±4.6	20.1±1.6	15.4±2.5	19.6±0.2	11.7±0.7	9.8±1.8
	OE ↓	10.9±0.8	19.6±0.1	22.0±4.6	19.8±1.6	14.3±2.4	18.9±0.1	10.6±0.4	9.0±1.6
iWildCam	AdaECE ↓	21.0±3.2	23.0±0.5	25.5±0.7	20.3±1.1	13.0±2.5	20.6±1.8	9.7±0.3	12.6±1.4
	OE ↓	14.3±3.6	16.8±0.2	20.4±0.8	15.5±0.2	8.21±1.4	15.4±1.2	5.2±0.2	7.9±1.0
FmoW	AdaECE ↓	20.0±9.9	20.9±8.6	41.7±0.1	22.4±0.9	9.73±0.1	21.7±0.2	10.5±0.2	10.6±0.1
	OE ↓	13.7±7.7	14.2±7.0	33.7±0.2	16.7±0.8	4.78±0.1	14.6±0.2	5.5±0.3	5.4±0.3

Table 6: We report additional OOD test scores (%) for reruns evaluated on both synthetic and wild benchmarks for *Peacock* and its components.

Dataset	ECE ↓	CE	MaxEnt	AdaFocal	RankMixup	CPC	Dual	ACLS	Peacock (Eq.)	Peacock (Impt.)
Amazon	Pre	7.0±0.5	5.0±0.6	6.7±1.0	43.1±0.1	7.4±0.5	5.8±2.3	42.0±0.1	6.5±3.2	5.0±1.1
	Post	11.6±0.5	7.6±0.4	10.0±0.3	41.0±0.2	3.4±1.2	5.0±1.5	26.3±0.1	5.6±1.7	4.8±0.2
	Avg.	9.3±0.5	6.3±0.2	8.4±0.2	42.0±0.3	5.4±0.3	5.4±0.2	34.2±0.3	6.1±0.3	4.9±0.3
	Temp.	1.50	1.25	1.25	1.25	1.25	1.25	1.25	1.25	1.25
CivilComments	Pre	10.4±0.4	4.8±0.2	7.8±1.7	11.4±0.5	2.4±0.4	4.2±0.1	11.1±0.1	2.1±0.8	4.2±1.0
	Post	6.3±0.9	8.2±0.7	11.6±0.2	11.0±0.1	2.2±0.3	7.5±0.4	6.7±0.1	5.7±0.7	7.5±0.5
	Avg.	8.4±0.6	6.5±0.2	9.7±1.5	11.2±0.5	2.3±0.4	5.9±0.1	8.9±0.2	3.9±0.9	5.8±0.9
	Temp.	2.00	1.25	1.25	1.25	1.25	1.25	1.25	1.25	1.25

Table 7: Vanilla temperature scaling results with temperatures obtained post-grid search for Wilds-Text datasets.

Algorithm	Acc (%)	ECE (%)	Speed (Sec)	w1	w2	w3	$\sum_t w_t = 1$
Equal-Importance	76.8±0.2	6.3±0.1	48.3±0.2	0.33	0.33	0.33	Yes
MTAN (Liu et al., 2019a)	76.5±0.1	6.5±0.1	48.5±0.2	0.33	0.33	0.33	Yes
CoVV (Groenendijk et al., 2021)	77.3±0.1	6.5±0.1	48.9±0.1	0.01	0.52	0.47	Yes
GradNorm (Chen et al., 2018)	75.7±0.6	6.8±0.4	79.1±0.1	2.99	0.00	0.00	No
MT-MOO (Sener & Koltun, 2018)	76.4±0.4	6.5±0.1	66.3±0.3	0.36	0.47	0.16	Yes
Weighted-Importance (Ours)	77.3±0.4	6.2±0.3	48.5±0.2	0.00	0.53	0.47	Yes

Table 8: Comparisons of different multi-objective optimization methods for *Peacock* evaluated on CIFAR10/CIFAR10-C using ResNet-18.

Algorithm (CIFAR10-C)	Acc (%)	ECE (%)	Speed (Sec)	w1	w2	w3	$\sum_t w_t = 1$
Equal-Importance	82.7±0.1	9.8±0.3	838±3	0.33	0.33	0.33	Yes
CoVV (Groenendijk et al., 2021)	83.1±0.2	8.5±0.3	827±5	0.02	0.22	0.76	Yes
GradNorm (Chen et al., 2018)	80.2±0.4	9.7±0.4	1347±3	3.00	0.00	0.00	No
MT-MOO (Sener & Koltun, 2018)	80.7±0.1	9.8±0.1	915±3	0.43	0.54	0.03	Yes
Weighted-Importance (Ours)	81.0±0.4	6.1±0.3	840±5	0.00	0.53	0.47	Yes

Table 9: Comparisons of different multi-objective methods for *Peacock* using SWINV2.

D.3 ABLATION STUDIES

To gain a better understanding of each component in *Peacock*, we provide an ablation study that removes each component from the full combination of *Peacock*. In Table 10, we show the respective ECE, OE and KSE scores of each combination evaluated on CIFAR/CIFAR-C. While each component generally helps improve calibration performance, we identify RankMixup and MaxEnt loss as two of the most critical building blocks of *Peacock*. Since the removal of either RankMixup or MaxEnt loss would cause a noticeable drop in calibration performance. Although ACLS independently delivers competitive performance, we find it to be the least impactful, since its removal leads to better calibration in *Peacock*. Therefore we propose the final version of equal and importance weighted forms *Peacock* to be without ACLS. Note that the experiments performed in this ablation study does not include temperature scaling. For e.g., removing RankMixup, would cause the highest ECE on CIFAR10/CIFAR10-C with 1.9% and 10.1% respectively. The lack of MaxEnt loss constraints delivers the worst result on CIFAR100/CIFAR100-C with 8.4% and 15.3%.

Algorithm (ID Performance)	(a) CIFAR10			(b) CIFAR100		
	ECE	NLL	KSE	ECE	NLL	KSE
Peacock w/o MaxE	1.2 \pm 0.1	271.8 \pm 2.9	1.1 \pm 0.1	8.4 \pm 0.1	316.8 \pm 1.7	8.4 \pm 0.1
Peacock w/o AdaFocal	1.7 \pm 0.1	289.5 \pm 4.2	1.8 \pm 0.1	6.1 \pm 0.7	314.8 \pm 1.0	6.1 \pm 0.7
Peacock w/o RankMixup	1.9 \pm 0.2	294.8 \pm 3.2	2.0 \pm 0.2	6.1 \pm 0.3	316.6 \pm 0.9	6.2 \pm 0.3
Peacock w/o CPC	1.6 \pm 0.2	284.6 \pm 6.3	1.9 \pm 0.2	6.0 \pm 0.7	317.8 \pm 1.9	6.0 \pm 0.7
Peacock w/o Dual	1.5 \pm 0.1	278.2 \pm 2.9	1.7 \pm 0.1	6.5 \pm 0.2	308.6 \pm 0.9	6.5 \pm 0.2
Peacock w/o ACLS	0.6 \pm 0.1	240.9 \pm 0.4	0.9 \pm 0.1	6.5 \pm 0.2	306.3 \pm 0.9	6.8 \pm 0.2

Algorithm (OOD Performance)	(a) CIFAR10-C			(b) CIFAR100-C		
	ECE	NLL	KSE	ECE	NLL	KSE
Peacock w/o MaxE	7.6 \pm 0.4	270.4 \pm 2.9	7.2 \pm 0.4	15.3 \pm 0.1	360.0 \pm 0.5	14.9 \pm 0.1
Peacock w/o AdaFocal	8.6 \pm 0.4	286.0 \pm 1.3	8.3 \pm 0.4	12.4 \pm 0.6	355.9 \pm 1.3	12.2 \pm 0.6
Peacock w/o RankMixup	10.1 \pm 0.4	296.8 \pm 2.8	9.9 \pm 0.5	12.7 \pm 0.4	358.4 \pm 0.4	12.4 \pm 0.4
Peacock w/o CPC	8.7 \pm 0.4	285.0 \pm 1.8	8.3 \pm 0.3	12.4 \pm 0.7	359.9 \pm 1.7	12.3 \pm 0.7
Peacock w/o Dual	8.0 \pm 0.1	278.7 \pm 0.8	7.7 \pm 0.1	12.5 \pm 0.2	353.2 \pm 1.1	12.3 \pm 0.2
Peacock w/o ACLS	6.5 \pm 0.2	245.6 \pm 0.4	6.3 \pm 0.1	11.6 \pm 0.3	358.3 \pm 0.5	11.7 \pm 0.5

Table 10: Component analysis of *Peacock* reveals the best performance when all algorithms except ACLS are combined.

E LIMITATIONS

Component Permutations In the case of *Peacock*, we featured a total of seven baselines which gives a total of $2^7 - 1$ permutations. While the primary focus of our paper is looking at whether different calibration algorithms can be successfully combined, we constrained *Peacock* to the seven featured algorithms so as to keep experiments manageable. We note that there are many potential algorithms in the calibration family that could become promising candidates (see Fig. 2a).

Modularity and Future Components To the best of our ability, we built *Peacock* based on the most relevant SOTA calibration components. For each algorithm, we closely referenced the source code provided by the respective authors. As we believe that *Peacock* will perform as well/better than the average of its components, we specifically built *Peacock* in a modular fashion allowing the easy integration of future methods.

1566 F REPRODUCIBILITY CHECKLIST

1567

1568 If needed, we provide the reproducibility checklist of this paper.

1569

1570 This paper:

1571

1572 • Includes a conceptual outline and/or pseudocode description of AI methods introduced
(yes)

1573

1574 • Clearly delineates statements that are opinions, hypothesis, and speculation from objective
facts and results (yes)

1575

1576 • Provides well marked pedagogical references for less-familiare readers to gain background
necessary to replicate the paper (yes)

1577

1578

1579 Does this paper make theoretical contributions? (yes)

1580

1581 If yes, please complete the list below.

1582

1583 • All assumptions and restrictions are stated clearly and formally. (yes)

1584

1585 • All novel claims are stated formally (e.g., in theorem statements). (yes)

1586

1587 • Proofs of all novel claims are included. (yes)

1588

1589 • Proof sketches or intuitions are given for complex and/or novel results. (yes)

1588

1589 • Appropriate citations to theoretical tools used are given. (yes)

1589

1590 • All theoretical claims are demonstrated empirically to hold. (yes)

1591

1592 • All experimental code used to eliminate or disprove claims is included. (yes)

1592

1593 Does this paper rely on one or more datasets? (yes)

1594

1595 If yes, please complete the list below.

1596

1597 • A motivation is given for why the experiments are conducted on the selected datasets (yes)

1598

1599 • All novel datasets introduced in this paper are included in a data appendix. (NA)

1600

1601 • All novel datasets introduced in this paper will be made publicly available upon publication
of the paper with a license that allows free usage for research purposes. (NA)

1601

1602 • All datasets drawn from the existing literature (potentially including authors' own previ-
ously published work) are accompanied by appropriate citations. (yes)

1603

1604 • All datasets drawn from the existing literature (potentially including authors' own previ-
ously published work) are publicly available. (yes)

1605

1606 • All datasets that are not publicly available are described in detail, with explanation why
publicly available alternatives are not scientifically satisficing. (NA)

1607

1608

1609 Does this paper include computational experiments? (yes)

1610

1611 If yes, please complete the list below.

1612

1613 • Any code required for pre-processing data is included in the appendix. (yes).

1614

1615 • All source code required for conducting and analyzing the experiments is included in a
code appendix. (yes)

1616

1617 • All source code required for conducting and analyzing the experiments will be made pub-
licly available upon publication of the paper with a license that allows free usage for re-
search purposes. (yes)

1618

1619 • All source code implementing new methods have comments detailing the implementation,
with references to the paper where each step comes from (yes)

- 1620 • If an algorithm depends on randomness, then the method used for setting seeds is described
1621 in a way sufficient to allow replication of results. (yes)
- 1622 • This paper specifies the computing infrastructure used for running experiments (hardware
1623 and software), including GPU/CPU models; amount of memory; operating system; names
1624 and versions of relevant software libraries and frameworks. (yes)
- 1625 • This paper formally describes evaluation metrics used and explains the motivation for
1626 choosing these metrics. (yes)
- 1627 • This paper states the number of algorithm runs used to compute each reported result. (yes)
- 1628 • Analysis of experiments goes beyond single-dimensional summaries of performance (e.g.,
1629 average; median) to include measures of variation, confidence, or other distributional in-
1630 formation. (yes)
- 1631 • The significance of any improvement or decrease in performance is judged using appropri-
1632 ate statistical tests (e.g., Wilcoxon signed-rank). (yes)
- 1633 • This paper lists all final (hyper-)parameters used for each model/algorithm in the paper's
1634 experiments. (yes)
- 1635 • This paper states the number and range of values tried per (hyper-) parameter during devel-
1636 opment of the paper, along with the criterion used for selecting the final parameter setting.
1637 (yes)
- 1638
- 1639
- 1640
- 1641
- 1642
- 1643
- 1644
- 1645
- 1646
- 1647
- 1648
- 1649
- 1650
- 1651
- 1652
- 1653
- 1654
- 1655
- 1656
- 1657
- 1658
- 1659
- 1660
- 1661
- 1662
- 1663
- 1664
- 1665
- 1666
- 1667
- 1668
- 1669
- 1670
- 1671
- 1672
- 1673