# ANALYTICAL SOFTMAX TEMPERATURE SETTING FROM FEATURE DIMENSIONS FOR MODEL- AND DOMAIN-ROBUST CLASSIFICATION *

**Tatsuhito Hasegawa**
Graduate School of Engineering,
University of Fukui,
Fukui
t-hase@u-fukui.ac.jp

**Shunsuke Sakai**
Graduate School of Engineering,
University of Fukui,
Fukui
mf240599@g.u-fukui.ac.jp

## ABSTRACT

In deep learning-based classification tasks, the softmax function's temperature parameter $T$ critically influences the output distribution and overall performance. This study presents a novel theoretical insight that the optimal temperature $T^*$ is uniquely determined by the dimensionality of the feature representations, thereby enabling training-free determination of $T^*$. Despite this theoretical grounding, empirical evidence reveals that $T^*$ fluctuates under practical conditions owing to variations in models, datasets, and other confounding factors. To address these influences, we propose and optimize a set of temperature determination coefficients that specify how $T^*$ should be adjusted based on the theoretical relationship to feature dimensionality. Additionally, we insert a batch normalization layer immediately before the output layer, effectively stabilizing the feature space. Building on these coefficients and a suite of large-scale experiments, we develop an empirical formula to estimate $T^*$ without additional training while also introducing a corrective scheme to refine $T^*$ based on the number of classes and task complexity. Our findings confirm that the derived temperature not only aligns with the proposed theoretical perspective but also generalizes effectively across diverse tasks, consistently enhancing classification performance and offering a practical, training-free solution for determining $T^*$.

***Keywords*** Temperature Determination Coefficients, Optimal Temperature Parameter, Softmax Function, Training-free hyperparameter selection

## 1 Introduction

Deep learning exhibits impressive performance in various areas, such as image recognition [1], object detection [2], image generation [3], and natural language processing [4]. In general, a deep learning model can be constructed by stacking layers, such as a convolutional layer, a fully connected layer, and a normalization layer. This model is then trained by minimizing the expected risk using a defined loss function. The cross-entropy loss displayed in Eq. (1) is typically used for standard classification:

$$\mathcal{L}(f(\boldsymbol{x};\boldsymbol{\theta}),\boldsymbol{y}) = -\frac{1}{n}\sum_{i=1}^{n}\sum_{j=1}^{c} y_{ij}\log f_j(\boldsymbol{x}_i;\boldsymbol{\theta}), \tag{1}$$

where $n$ and $c$ represent the number of samples and categories, respectively; $\boldsymbol{\theta}$ denotes the parameters of the deep neural network; $f_j(\boldsymbol{x}_i;\boldsymbol{\theta})$ denotes the $j$-th output of the model corresponding to input $\boldsymbol{x}_i$; and $y_{ij} \in \{0, 1\}$ indicates whether the $i$-th sample belongs to category $j$. Because the cross-entropy loss takes a true one-hot distribution $\boldsymbol{y}$ and an estimated distribution, the output of the model $\hat{\boldsymbol{y}}$ is normalized by the softmax function, as expressed in Eq. (2):

$$\text{softmax}(\hat{y}_j, \hat{\boldsymbol{y}}) = \frac{\exp(\frac{\hat{y}_j}{T})}{\sum_{\hat{y}_k \in \hat{\boldsymbol{y}}} \exp(\frac{\hat{y}_k}{T})}, \tag{2}$$

where $T$ is a hyperparameter of the softmax function, known as the temperature. In Fig. 1, we illustrate the effect of different values of $T$ when the model's output $\hat{\boldsymbol{y}}$ is fixed. Intuitively, the temperature $T$ controls the uncertainty of the estimated distribution. As $T$ increases, the estimated distribution approaches a uniform distribution; conversely, as $T$ decreases, the estimated distribution becomes deterministic. In many settings, the temperature $T$ is set to a default value of 1.0. Examples of applications involving $T$ include knowledge distillation [5], which transfers knowledge from a larger model to a smaller model, and temperature scaling [6], which adjusts the model's confidence. While some applications control the temperature to meet specific requirements, the effect of temperature on the model's learning process remains poorly understood. However, the temperature parameter $T$ used during the training process is a crucial factor that significantly impacts the model's final generalization performance. T adjusts the sharpness of the probability distribution output by the softmax function, which effectively scales the gradients calculated through the cross-entropy loss. This gradient scaling directly influences the learning dynamics, convergence stability, and ultimately, how well the model performs on unseen data (i.e., generalization performance). Furthermore, an appropriate value of $T$ can act as a form of regularization, potentially preventing the model from overfitting to the training data and helping it learn more generalizable representations. Therefore, **selecting a suitable $T$ for the task and model is essential not only for adjusting the model's output distribution but also for optimizing the training process itself to build models with better generalization capabilities**. Particularly in many practical scenarios where high generalization performance on unseen data is critically important, such as large-scale image classification [1], natural language understanding [4], or time-series forecasting in changing environments [7], optimizing the temperature parameter during training is crucial for enhancing model reliability and utility.
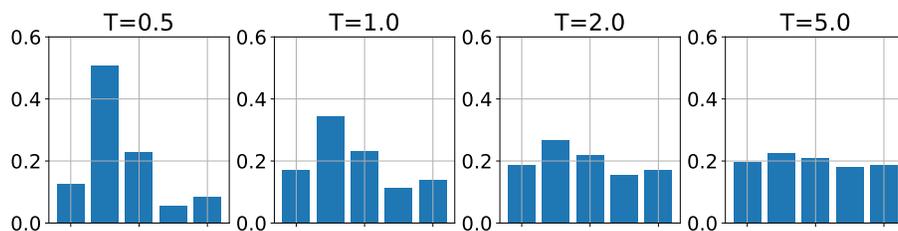


Figure 1: Sample outputs of the softmax function at various temperatures for a fixed model output $\hat{\boldsymbol{y}}$.

Some studies [8, 9] investigated the effect of changes in temperature parameters on model training. Agarwala *et al.* [8] tackled image classification tasks by introducing the inverse temperature $\beta$ as a hyperparameter and, through extensive experimentation on various tasks, demonstrated that tuning $\beta$ can significantly improve model performance. While they established the importance of $\beta$ for optimizing generalization, no universal formula or method was put forward to pinpoint the best $\beta$ across disparate model architectures and datasets. Concurrently, Hasegawa [9] explored a possible correlation between the optimal temperature $T$ and the dimensionality $M$ of the feature map in a 1D CNN-based human activity recognition scenario and showed a trend in which $T$ grows with $M$, suggesting that these two factors may be functionally linked. However, variations in tasks and model architectures introduced inconsistencies, rendering it difficult to derive a single, general rule for determining $T$. Taken together, these findings underscore both the necessity and the complexity of fine-tuning the temperature parameter. Despite these advances, the community still lacks a systematic, generalizable method to determine either $T$ or $\beta$ robustly in the face of architectural and task-related differences.

In this study, we propose a novel method for estimating the optimal temperature parameter $T^*$ in a *training-free* manner that is robust to variations in both model architectures and tasks. By eliminating the need for extensive tuning, our approach not only reduces the computational burden but also achieves superior performance compared to the commonly used default $T = 1$. Because minimizing hyperparameters is critical for the practical deployment of deep learning models, this study represents a significant step toward more efficient and scalable deep learning solutions. Since the optimal value of the temperature $T$, a hyperparameter directly involved in the training process, typically depends heavily on the dataset and model architecture [8, 10], tuning it conventionally requires computationally expensive trial-and-error with multiple full model training runs. Our method is the first to provide a training-free, closed-form estimator for the optimal temperature $T^*$, leveraging feature dimensionality and task characteristics. By eliminating costly searches, it offers significant practical value—substantially saving computational resources and development time—while also introducing a principled theoretical framework for temperature determination in deep learning.

The main contributions of this study are as follows:

- **A theoretical framework linking $T^*$ and feature dimensionality.** We establish a theoretical relationship between the optimal temperature parameter $T^*$ and the feature dimensionality $M$. In addition, we propose to represent $T^*$ using a set of *temperature determination coefficients* $\alpha, \beta, \gamma, \delta$ as

$$T^* = \alpha\sqrt{M} \ + \ \beta \ + \ \gamma \log(csg) \ + \ \delta \log(cn),$$

  where $\alpha, \beta, \gamma, \delta$ are experimentally optimized constants, and $csg$, $cn$ denote the cumulative spectral gradient (CSG) [11] and the number of classes (CN), respectively. Furthermore, we show that inserting a batch normalization (BN) layer [12] immediately before the output layer effectively stabilizes this relationship, making $T^*$ estimation less sensitive to changes in model architecture and dataset characteristics.

- **Empirical analysis of BN insertion and its performance impact.** Through extensive experiments, we demonstrate that while placing a BN layer right before the output layer may cause performance degradation under the conventional setting $T = 1$, it can improve performance when the temperature $T^*$ is set appropriately. This finding emphasizes BN's dual role of normalizing the feature space and facilitating more precise temperature selection, which in turn leads to improved generalization.

- **An empirical formula for training-free estimation of $T^*$.** By leveraging large-scale experiments, we optimize the *temperature determination coefficients* $\alpha, \beta$ to propose a closed-form solution for $\hat{T}^*$ that requires no additional training overhead:

$$\hat{T}^* = \mathrm{clip}\big(0.7239\sqrt{M} - 4.706, \ \epsilon, \ 512\big),$$

  which consistently outperforms the conventional default $T = 1$. This formula serves as a reliable starting point for determining $T^*$, eliminating the need for time-consuming hyperparameter searches.

- **Task-aware refinement using temperature determination coefficients.** We further refine $\hat{T}^*$ by incorporating two task-dependent terms: the number of classes ($cn$) and the cumulative spectral gradient ($csg$). Specifically, we again optimize the temperature determination coefficients $\alpha, \beta, \gamma, \delta$ using large-scale experimental results to obtain:

$$\hat{T}^*_{csgcn} = \mathrm{clip}\Big(0.3191\sqrt{M} + 20.74 + 3.746\log(csg) - 7.380\log(cn), \ \epsilon, \ 512\Big).$$

  By accounting for both task complexity and the number of classes, this refined scheme offers robust and generalizable performance gains across a wide range of deep learning applications.

## 2 Related Works

### 2.1 Temperature parameter

Several studies have utilized the temperature of the softmax function, with knowledge distillation (KD) [5, 13, 14, 15] and temperature scaling [6, 16] serving as representative examples. Knowledge distillation is a method for transferring knowledge from a teacher model to a student model. This method aims to improve performance and reduce the number of model parameters. In conventional KD, it is well known that setting the softmax temperature to approximately $T = 4$ during distillation improves performance. Liu *et al.* [13] proposed Meta KD, which optimizes $T$ via meta-gradients computed on a validation set. Curriculum Temperature KD (CTKD) [14] introduces an adversarial learning scheme that directly learns the temperature parameter itself. Sun *et al.* [15] achieve a dynamic temperature by applying logit standardization independently to both the teacher and the student models. Furthermore, temperature scaling [6] reduces the model's prediction bias by correcting temperature parameters. Balanya *et al.* [16] utilized adaptive temperature in the temperature scaling. In contrastive learning, the temperature parameter is also adjusted to modulate the influence of negative samples [17].

In addition to these approaches, Agarwala *et al.* [8] presented a temperature check method examining the effect of the inverse temperature $\beta = 1/T$ on model generalization, highlighting the importance of temperature adjustment. The authors tuned the inverse temperature during standard model training within the range $\beta \in [10^{-2}, 10^1]$. Xuan *et al.* [10] also highlighted the importance of tuning the temperature parameter in classification tasks. Furthermore, relaxed softmax [18] also proposes a tuning strategy for inverse temperature. In fields other than classification, temperature parameters have attracted attention in deep metric learning. Xu [19] proposed a heated-up strategy that involves training the model with increasing temperature.

However, none of the above studies discusses the relationship between the dimensionality of the feature maps $M$ and the optimal temperature $T^*$. While these studies note the importance of tuning the temperature parameter, they do not

sufficiently address the question of whether the optimal temperature parameter can be estimated without training the model.

Building upon these general insights into temperature tuning, a study in sensor-based human activity recognition has explored the relationship between feature dimensionality and temperature [9]. Denoting the number of dimensions of the output of the feature extractor (e.g., ConvNets) as $M$, the experimental results indicate a potential relationship between the optimal temperature $T^*$ and $M$. However, there were several limitations:

1. The experiments were conducted only on one-dimensional sensor-based human activity recognition datasets.
2. The coefficients determining the optimal temperature parameter vary across different datasets and the encoder's model architectures.
3. The effectiveness has not been evaluated using a deeper model.

For these limitations, the related study [9] discusses a method for mitigating limitation (2) by inserting layer normalization (LN) [20] before applying the softmax function. This method is inspired by the observation that an output distribution with optimal temperature $T^*$ follows a specific distribution. By normalizing the outputs and dynamically adjusting the temperature parameter through the trainable parameters in LN, this method can be regarded as an advanced extension of the relaxed softmax approach [18]. However, it does not account for the potential relationship between $M$ and $T^*$. In this paper, we further explore the relationship between $M$ and $T^*$.

## 2.2 Scaled dot-product attention

The scaled dot-product attention, introduced by Vaswani *et al.* [21], focuses on the temperature parameter and the number of input dimensions of the softmax function. This operation can be computed as follows:

$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V, \tag{3}$$

where $Q$ and $K^T$ denote the feature map $\boldsymbol{z}$ and parameters $\boldsymbol{w}_j$, respectively. Here, the problem of increasing dispersion occurs as the number of dimensions increases, as discuss in the next section. The authors addressed this problem by dividing by $\sqrt{d_k}$, the square root of the number of dimensions of $Q$ and $K$. This is equivalent to the softmax function with $T = \sqrt{d_k}$. The authors introduced this temperature adjusting to maintain an appropriate gradient scale. However, they did not provide a detailed discussion of this temperature adjustment.

## 2.3 Label smoothing

Label smoothing [22] is a regularization method that has a role similar to that of temperature scaling. As displayed in Eq. (1), a typical cross-entropy loss function employs ground-truth labels $\boldsymbol{y} \in \{0, 1\}^c$, which are encoded in a one-hot manner. In contrast, label smoothing represents ground-truth labels $\boldsymbol{y}^{LS} \in \mathbb{R}^c$ in a soft manner, as expressed in Eq. (4):

$$y_i^{LS} = y_i(1 - \epsilon) + \epsilon/c. \tag{4}$$

Label smoothing minimizes interclass variance, maximizes intraclass variance, and improves the model's generalization ability [23]. In addition to a method for statically constructing soft labels from a uniform distribution, a dynamic label construction method has also been proposed [24].

Temperature adjusting and label smoothing are similar in that they focus on the distribution of negative samples. However, whereas temperature adjusting controls the estimated distribution $\hat{\boldsymbol{y}}$, label smoothing transforms the ground-truth distribution $\boldsymbol{y}$. In Fig. 2, we visualize a simulation of the response of the cross-entropy loss as the logits $\hat{y}_i$ are varied. In the case without label smoothing, we observe that the loss monotonically decreases, with temperature regulating extreme changes in the loss. In the case of label smoothing, we observe that label smoothing penalizes large logits $\hat{y}_i$ in response to $\epsilon$. As illustrated in the figure, the effects of both regularizations are observed in the case involving label smoothing ($\epsilon = 0.1$).

Based on the above observations, we conclude that label smoothing and temperature adjusting serve distinct roles in model training. Specifically, we focus on the relationship between the dimensionality of the feature map $M$ and the optimal temperature $T^*$.

## 2.4 Position of this study

In this section, we clarify the position of our study. Although knowledge distillation [5] and temperature scaling [6] have demonstrated the effectiveness of tuning the temperature parameter, both methods primarily focus on distillation
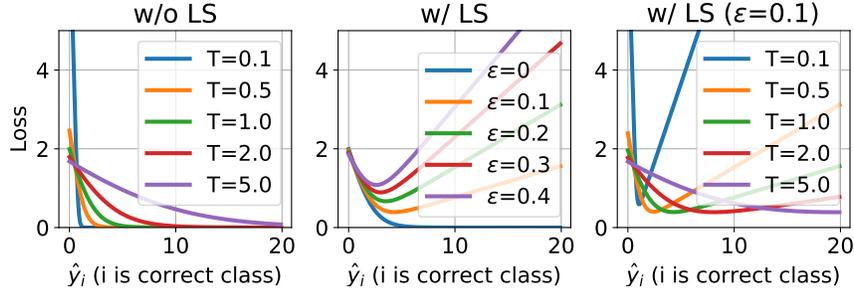
Figure 2: Effect of temperature variation and label smoothing on cross-entropy loss.

and calibration rather than classification. Scaled dot-product attention [21], on the other hand, replaces the temperature with the square root of the feature dimensionality $d_k$ in the softmax function. However, it is treated as part of the model architecture and does not involve deriving an explicit optimal temperature for classification. Furthermore, although relaxed softmax [18] and inserting LN [9] dynamically adjusts the temperature parameter during training, a theoretical framework linking the feature dimensionality $M$ to the optimal temperature $T^*$ has not been established. As such, existing research has not theoretically established the optimal temperature $T^*$ for improving generalization performance in classification tasks.

To address this gap, we propose a new approach that does not require additional training and incorporates the effect of $M$ into a closed-form estimation of the optimal temperature. As summarized in Table 1, this study formalizes and exploits the previously underexplored relationship between feature dimensionality and the temperature parameter, setting it apart from conventional methods. Our goal is to achieve consistent performance improvements across various model architectures and application tasks by offering a theoretically grounded and practically efficient temperature estimation procedure. In particular, while temperature check [8] determines $T^*$ through costly hyperparameter searches involving extensive model retraining, our proposed method is novel in that it uniquely identifies $T^*$ in a training-free manner. Although Hasegawa *et al.* [9] also noted a relationship between $M$ and $T$, they did not derive a closed-form solution. By contrast, our approach defines an explicit closed-form solution and further incorporates corrections for task difficulty and class count, representing a novel contribution.

Table 1: Characteristics of our method compared with related works.

| Method | Task | Tuning | Consider $M$ | Closed-form |
|---|---|---|---|---|
| KD [5] | Distillation | Required | No | No |
| Meta KD [13] | Distillation | Dynamic | No | No |
| CTKD [14] | Distillation | Dynamic | No | No |
| Logit standardization [15] | Distillation | Dynamic | No | No |
| Temperature Scaling [6] | Calibration | Required | No | No |
| Temperature Check [8, 10] | Classification | Required | No | No |
| Relaxed Softmax [18] | Classification | Dynamic | No | No |
| Scaled Dot-Product Attn [21] | Architecture | Unnecessary | Yes | Yes |
| Label Smoothing [22] | Classification | Required | No | No |
| Insert LN [9] | Classification | Dynamic | No | No |
| Ours | Classification | **Unnecessary** | **Yes** | **Yes** |

## 3 Training-free temperature determination

As demonstrated in related works, tuning the temperature parameter of the softmax function is critical for improving performance in classification tasks [8, 9]; however, its optimal value has traditionally been chosen empirically—either fixed at 1 or selected via exploratory methods such as grid search. In this study, we introduce a training-free temperature determination approach that identifies the optimal temperature without any hyperparameter search. Our method is exceptionally easy to integrate, requires no modifications to existing training pipelines, and is broadly applicable to any deep learning model employing softmax cross-entropy loss. Moreover, it rests on a transparent theoretical framework, making its behavior both predictable and interpretable.

Implementing our approach entails only the following simple steps:

1. Inserting BN just before the output layer as illustrated in Fig. 3.
2. Determining the temperature by following equation:

$$T^* = \underbrace{\alpha \sqrt{M}}_{\text{Model correction}} + \underbrace{\beta}_{\text{Bias term}} + \underbrace{\gamma \log(\text{csg}) + \delta \log(\text{cn})}_{\text{Task correction}}, \tag{5}$$

where $\alpha, \beta, \gamma, \delta$ are extended-version of temperature determination coefficients, and $csg$, $cn$ denote the CSG [11] and the number of classes, respectively. The first term compensates for changes in the feature-map dimensionality $M$, based on the hypothesis that $M$ exerts a dominant influence on the determination of $T^*$ (see Section 3.1). Furthermore, under the hypothesis that this dimensionality correction alone cannot fully remove the effects of differing model architectures (e.g., ResNet, ViT, etc), we insert a batch-normalization layer immediately before the output layer (see Section 3.2). The second term serves as a bias to adjust the overall baseline. The third and fourth terms correct for differences in task characteristics (see Section 3.3).
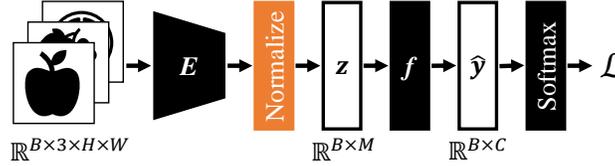


Figure 3: Outline of the flow of the general deep neural network models inserted a normalization layer.

## 3.1 Model correction by $\alpha\sqrt{M}$

$T^*$ has been shown to depend on the feature-map dimensionality $M$ in related work [9]. In this section, we theoretically motivate the need to introduce the model-correction term (the first term) in Eq. (5).

In general, a linear classifier computes the logits from the output $\boldsymbol{z} \in \mathbb{R}^M$ of the feature extractor $E$ as follows:

$$\hat{y}_j = \boldsymbol{W}_j \boldsymbol{z}^T + b_j, \tag{6}$$

where $\boldsymbol{W}_j = [w_{j1}, w_{j2}, \ldots, w_{jM}]$ is an $M$-dimensional weight vector corresponding to category $j$.

At a given training step, the expectation and variance of $\hat{y}_j$ are given by:

$$\mathbb{E}[\hat{y}_j] = \mathbb{E}\left[\sum_{k=1}^M w_{jk} z_k + b_j\right] = \sum_{k=1}^M w_{jk} \mathbb{E}[z_k] + b_j, \tag{7}$$

$$\mathbb{V}[\hat{y}_j] = \mathbb{V}\left[\sum_{k=1}^M w_{jk} z_k + b_j\right] = \mathbb{V}\left[\sum_{k=1}^M w_{jk} z_k\right]. \tag{8}$$

Here, $w$ is treated as a constant, and $z$ is a random variable that is generally neither independent nor normalized. In contrast, at the initial stage of training, $w$ and $z$ can often be approximated as independent and identically distributed. Under this assumption, the expectation and variance can be simplified to:

$$\mathbb{E}[\hat{y}_j] = M \mathbb{E}[w_j z] + b_j, \tag{9}$$
$$\mathbb{V}[\hat{y}_j] = M \mathbb{V}[w_j z]. \tag{10}$$

Because this variance grows proportionally to $M$, increasing model depth can lead to exploding or vanishing gradients. To address this issue, LeCun's initialization [25] sets the initial parameters according to $\mathbb{V}[w] = 1/M$, an approach that was subsequently extended in the Xavier [26] and He [27] initialization methods.

Focusing on the effect at the output layer, $\hat{y}_j$ is passed into the softmax function in Eq. (2). In other words, if $\hat{y}_j$ **depends on** $M$, it **implicitly modifies** the temperature parameter $T$ as a function of $M$. Consequently, this suggests that $T$ should be **chosen appropriately based on** $M$. In this study, we introduce **temperature determination coefficients** ($\alpha$ and $\beta$) to determine the optimal temperature $T^*$ as a function of $M$ as follows:

$$T^* = \alpha\sqrt{M} + \beta. \tag{11}$$

The variance of $\hat{y}_j/T^*$ within the softmax function is expressed by the following equation.

$$\mathbb{V}[\hat{y}_j/T^*] = \frac{1}{(\alpha\sqrt{M}+\beta)^2}\mathbb{V}\Big[\sum_{k=1}^{M} w_{jk}z_k\Big]. \tag{12}$$

When $\alpha = 1.0$ and $\beta = 0.0$, the variance becomes $1/M$, thereby suppressing the influence of $M$ on the softmax function. On the other hand, while $\mathbb{V}[\hat{y}_j/T^*]$ is not affected by $M$, it cannot be guaranteed to represent the optimal distribution for training when using the softmax cross-entropy loss. Therefore, the optimal $T^*$ for training is determined by adjusting the temperature determination coefficients.

### 3.2 Inseting BN just before the output layer

#### 3.2.1 Implementation

As illustrated in Fig. 3, The model architecture is simple, inserting the normalization layer just before the output layer. Considering image recognition as an example, the standard model takes an input image $x \in \mathbb{R}^{B \times 3 \times H \times W}$, extracts the feature maps $z \in \mathbb{R}^{B \times M}$ using the feature extractor $E$, and ultimately obtains the output $y \in \mathbb{R}^{B \times C}$ using the classifier $f$. In this case, we apply a global average pooling (GAP) layer to the output of the feature extractor. We can use any model architecture as a feature extractor, such as ConvNets and vision transformers (ViTs). For simplicity, we use a single fully connected layer for the classifier $f$.

#### 3.2.2 Hypothesis

In the related study [9], comprehensive validation experiments demonstrated that when the model architecture or task differs, $T^*$ behaves as a function of $M$, though its trend undergoes slight variations (i.e., the temperature determination coefficients fluctuate). We hypothesize that this may be due to variations in $z_k$ in Eq. (12) depending on the model or task.

In this study, we consider normalizing $z$ to enhance the robustness of temperature determination coefficients against variations in model architecture and task.

First, by transforming Eq. (8), the following equation can be derived.

$$\mathbb{V}\Big[\sum_{k=1}^{M} w_{jk}z_k\Big] = \mathbb{E}\Big[\Big\{\sum_{k=1}^{M} w_{jk}z_k\Big\}^2\Big] - \mathbb{E}\Big[\Big\{\sum_{k=1}^{M} w_{jk}z_k\Big\}\Big]^2$$
$$= \sum_{k=1}^{M}\sum_{l=1}^{M} w_{jk}w_{jl}\mathbb{E}[z_k z_l] - \sum_{k=1}^{M}\sum_{l=1}^{M} w_{jk}w_{jl}\mathbb{E}[z_k]\mathbb{E}[z_l]. \tag{13}$$

Furthermore, if we divide the first term into diagonalized and non-diagonalized terms,

$$\mathbb{V}\Big[\sum_{k=1}^{M} w_{jk}z_k\Big] = \sum_{k=1}^{M} w_{jk}^2\mathbb{E}[z_k^2] + \sum_{\substack{k,l=1\\k\neq l}}^{M} w_{jk}w_{jl}\mathbb{E}[z_k z_l] - \sum_{k=1}^{M}\sum_{l=1}^{M} w_{jk}w_{jl}\mathbb{E}[z_k]\mathbb{E}[z_l]. \tag{14}$$

Let us assume that $z$ has been normalized by BN to have a mean of 0 and a variance of 1, which implies that $\mathbb{E}[z] = 0$ and $\mathbb{E}[z^2] = 1$. Therefore, Eq. (14) can be transformed as follows:

$$\mathbb{V}\Big[\sum_{k=1}^{M} w_{jk}z_k\Big] = \sum_{k=1}^{M} w_{jk}^2 + \sum_{k\neq l} w_{jk}w_{jl}\mathbb{E}[z_k z_l]. \tag{15}$$

Comparing Eqs. (14) and (15), it can be observed that the influence of $z$ disappears from the first term of Eq. (14), and the third term becomes zero by introducing the normalization. As a result, it can be stated that the impact of $z$ on the variance has been significantly reduced.

Eq. (15) indicates that the variance is affected by the covariance of $z$ without being influenced by $z$ itself. However, the effect of the non-diagonal elements is generally not significant because correlations between feature maps are low when $M$ is sufficiently large. Based on the above, we hypothesize that we can mitigate the effects of $\hat{y}_j$ and $z$ by applying normalization to $z$.

7

### 3.3 Task correction by $\gamma \log(\mathbf{csg}) + \delta \log(\mathbf{cn})$

In Eq. (15), it was demonstrated that the insertion of BN can reduce the influence of $z$ on the variance. However, while this influence is expected to be relatively minor, the second term of Eq. (15) still retains some dependency on $z$. Therefore, we hypothesize that the influence appearing in the second term is affected by factors such as task difficulty and the number of output classes. To address this, we aim to reduce the impact by extending the temperature determination coefficients, as shown in Eq. (5).

Note that the second term in Eq. (15) merely indicates that some inter–feature covariance may persist after BN, but it does not guarantee a one–to–one correspondence between that covariance and task difficulty. Indeed, when the latent representations are sufficiently separated, $\mathbb{E}[z_k z_l] \approx 0$ can still hold even for hard tasks. To capture both cases in a unified way, we incorporate two complementary correction viewpoints:

(a) *Representation-space correction:* Empirically, tasks with highly similar classes exhibit larger inter–feature correlations. We therefore adopt the CSG as a proxy for this effect (see Appendix B.1 for a detailed derivation).

(b) *Softmax-gradient correction:* As the number of classes $C$ increases, the expected true-class probability approaches $1/C$, diluting the cross-entropy gradient $\partial L/\partial z_i$. This gradient dilution can be compensated by reducing the temperature $T$ (see Appendix B.2 for a detailed derivation).

Combining these two viewpoints yields the final temperature correction rule

$$\gamma \log(\mathbf{csg}) \;+\; \delta \log(\mathbf{cn}),$$

where $\gamma > 0$ corrects for the increase in inter–feature correlation (i.e. task difficulty) and $\delta < 0$ offsets the softmax dilution caused by larger class counts.

## 4 Effect of BN insertion

In the subsequent sections, we first examine the impact of BN insertion in Section 4, then optimize the temperature determination coefficients in Section 5, and finally verify the effectiveness of the newly derived temperature determination method in Section 6.

In these experiments, we investigate the effect of inserting BN in image classification tasks. The research questions are as follows:

1. How does the insertion of BN just before the output layer affect estimation accuracy?
2. Does the insertion of BN just before the output layer stabilize the temperature determination coefficients?

### 4.1 Experimental settings

We utilized the CIFAR10/100 [28], STL10 [29], and Tiny ImageNet (Tiny IN) [30] datasets and employed VGG9 [31], ResNet10 [32], and PyramidNet10 [33] as model architectures. Due to resource constraints and the need for numerous trials, we selected relatively small datasets and shallow model architectures.

Our experiments utilized various models to equip the classifier $f$ with one fully connected layer. For all models, we initialized the weights of the feature extractor with the initialization proposed by He *et al.* [27] and the weights of the classifier with a uniform distribution $U(-1/\sqrt{M}, 1/\sqrt{M})$. We did not use bias terms in the convolutional layer. For model training, we used stochastic gradient descent (SGD) with a momentum of 0.9 and a weight decay of 0.0005. We set the initial learning rate to 0.1 and employed cosine annealing [34] as the learning rate scheduler. The number of epochs and batch size were set to 200 and 128, respectively. We applied standardization as preprocessing and employed cropping, horizontal flip, and rotation at random.

We increased or reduced $M$ by uniformly adjusting the number of filters in the model. For example, standard ResNet includes [64, 128, 256, 512] filters for each block. We reduced $M$ by dividing all filters by a factor of $n$. We report the median accuracy for multiple trials using different random seeds.

### 4.2 Effects of temperature in image recognition

In Fig. 4, we present the evaluation results of VGG9 ($M = 512$) on CIFAR-10. The performance improves as $T$ increases up to $T = 128$, after which it declines rapidly ($T \geq 256$ is outside the drawing range). As displayed in Fig. 1, lower values of $T$ lead to sharper $\hat{\mathbf{y}}$. Setting $T$ too low can cause the model to rapidly converge to local minima.

Conversely, higher values of $T$ lead to flatter $\hat{\boldsymbol{y}}$, allowing the model to explore a wider parameter space. Therefore, setting $T$ to an excessively high value may prevent the model from converging.
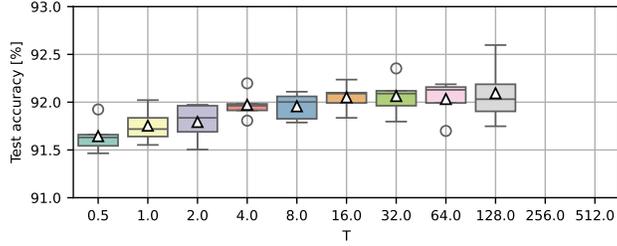


Figure 4: Change in accuracy relative to $T$ in the CIFAR10 environment using VGG9 with $M = 512$. The results for $T = 256$ and $T = 512$ are below the drawing range.

In Fig. 5, we present results across various model architectures and datasets with varying temperature parameters $T$. To standardize the CN, we extracted a subset of CIFAR100 containing superclasses at even positions (10 classes, 25,000 samples, denoted as CIFAR100@10). In each subfigure, the $x$-axis represents $T$, the $y$-axis represents $M$ for each model architecture, and the intensity represents the median test accuracy. A solid line denotes the best results, while a dashed line denotes the second-best results. From Fig. 5, we observe the trend that $\hat{T}^*$ increases as $M$ increases. However, as reported in the related study [9], this trend varies across datasets and model architectures.
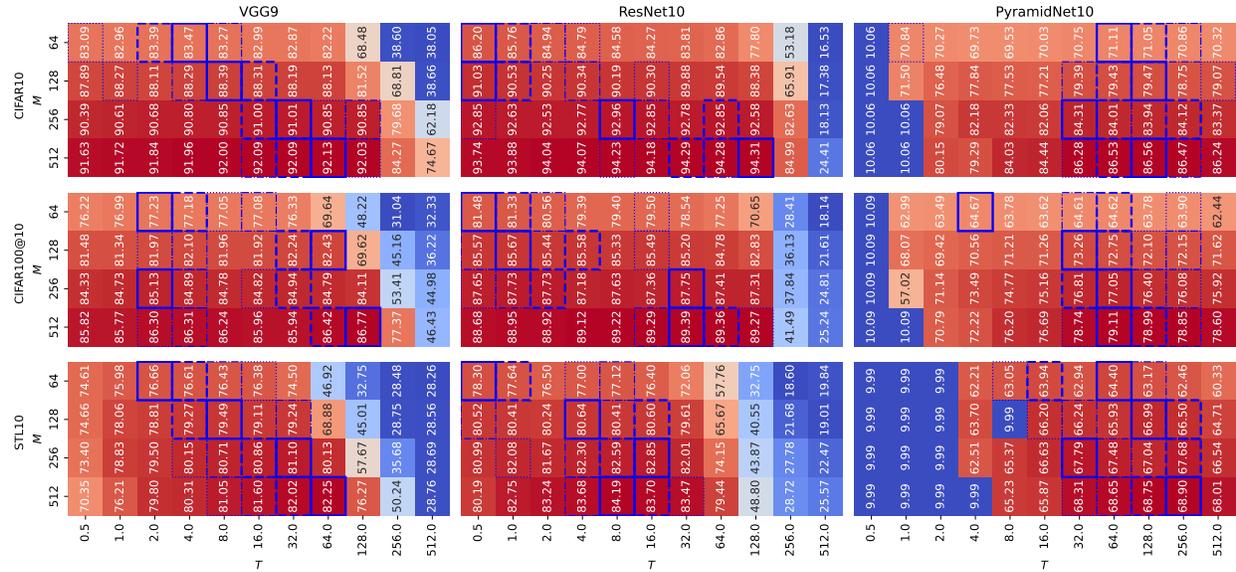


Figure 5: Test accuracies [%] for each temperature parameter in various scenarios (without insertion of the normalization layer). Only CIFAR-100@10 has 10 classes extracted from its superclasses (using only even class numbers) to standardize the number of output units to 10.

## 4.3 Effect of BN insertion

In Fig. 6, we present the results of inserting BN just before the output layer. This figure can be interpreted in the same way as Fig. 5. The trend of the optimal temperature parameter appears more robust than that displayed in Fig. 5. Notably, at the point where the performance declines in the region of high values of $T$, BN insertion causes the trend to be similar across different model architectures.

Next, we analyze the effect of BN insertion on estimation accuracy. In Fig. 7, we illustrate the differences in estimation accuracy with and without BN insertion (as $\Delta$ acc.). Positive values indicate that BN insertion improves performance, while negative values indicate that it reduces performance. As shown by the solid line, inserting BN at $T = 1$ shows that
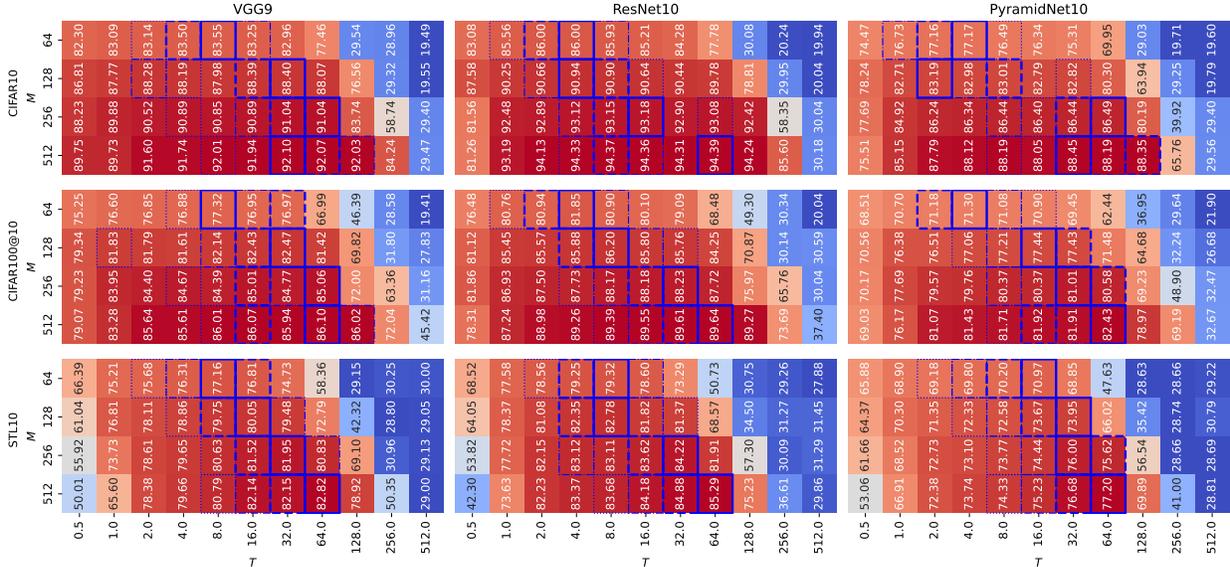
9

Figure 6: Test accuracies [%] for each temperature parameter in various scenarios (with BN insertion). Only CIFAR-100@10 has 10 classes extracted from its superclasses (using only even class numbers) to standardize the number of output units to 10.

the accuracy decreases with increasing M. However, by combining $T = \hat{T}^*$ and BN insertion, the accuracy improves in many cases. The performance improvement is particularly notable for a relatively high-$M$ and difficult dataset (STL-10).



Figure 7: Accuracy difference [%] between scenarios with and without BN insertion. Each value is calculated by subtracting the accuracy without BN from the accuracy with BN.

Based on the experimental results, inserting BN just before the output layer can provide a more robust estimation of the temperature determination coefficients across different datasets and model architectures. We infer the reason many conventional methods do not insert BN just before the output layer is that, with the default setting of $T = 1$, accuracy decreases due to BN insertion. However, it is possible to achieve a more accurate prediction once the temperature parameter $T$ is set appropriately. These results highlight the importance of optimizing the temperature parameter $T$.

10

# 5 Optimization of Temperature Determination Coefficients

In this section, we estimate the temperature determination coefficients through optimization. We experimentally demonstrate that (i) we can stabilize the estimation of the temperature determination coefficients by BN insertion and (ii) we can improve the estimation accuracy by optimizing the temperature parameter $T$. Based on these two results, we further introduce (iii) a correction of temperature using CSG and CN.

## 5.1 Estimation of performance for each $T$ and optimization of $T$

Since the optimal temperature $T^*$ under given conditions is unknown, we estimate the temperature determination coefficients $\alpha$ and $\beta$ as follows:

(i) We estimate the test performance at unknown $T$ by linear interpolation (Fig. 8).

(ii) We determine the optimal temperature using $\hat{T}^* = \alpha\sqrt{M} + \beta$.

(iii) By maximizing $\sum_{c \in C} \text{interpolation}_c(\hat{T}^*)$, we estimate the temperature determination coefficients $\alpha$ and $\beta$.
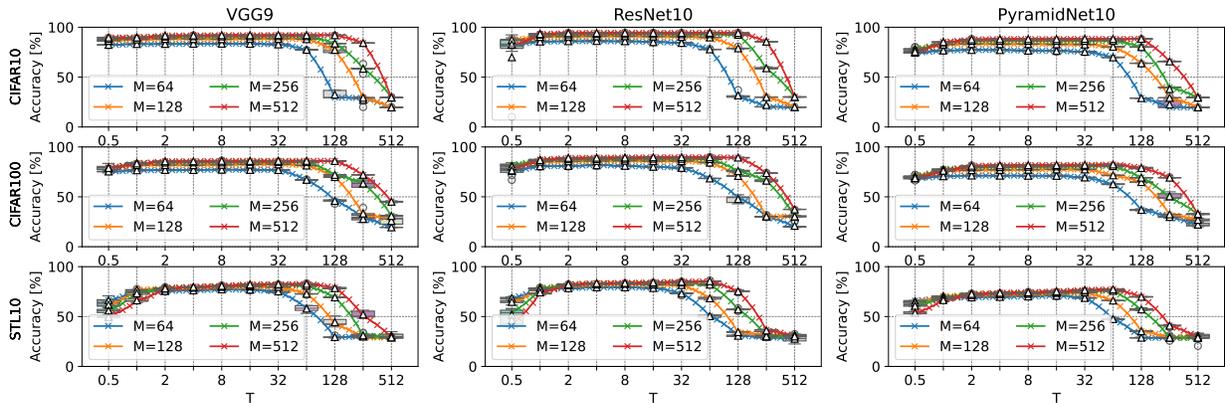


Figure 8: Estimation of test accuracy [%] for each temperature parameter (with BN insertion). The boxplots represent the observed values, and triangles indicate the median values. The solid line plots, marked with $\times$, represent the estimated values obtained through linear interpolation of the observed mean values.

In step (i), as illustrated in Fig. 8, we estimate the test performance at unknown $T$ by linear interpolation based on the results in Fig. 6. We present the results of Fig. 6 using a boxplot and draw a solid line indicating the test performance at unknown $T$ that is linearly interpolated. As the test performance does not significantly change with slight changes in $T$, we consider linear interpolation to be sufficient. This allows us to estimate the test performance in the range $0 < T \leq 512$. In step (ii), as mentioned in Section 3.2.2, we determine the optimal temperature $T^*$ by $\hat{T}^* = \alpha\sqrt{M} + \beta$. In step (iii), we estimate the temperature determination coefficients by maximizing the sum of the estimated test accuracy under various conditions $C$. We denote the estimated test accuracy for a given condition $C$ and temperature parameter $T$ as $\text{interpolation}_c(T)$. For example, we can estimate the optimal test accuracy of VGG9 on CIFAR-10 by setting the condition $C$ as $M = [64, 128, 256, 512]$. We employ differential evolution for optimization [35].

Table 2 displays the temperature determination coefficients estimated in each environment in Fig. 8 along with the test accuracy using these coefficients. We observe that, on average, $\alpha = 4$ and $\beta = -25$ with and without BN. However, the standard deviation of the estimated coefficient $\alpha$ is 2.17 without BN, which decreases to 1.86 with BN. Therefore, the stability of the coefficients improves with BN insertion. From the perspective of test performance, the test accuracy ($\mu_{acc}$) of the optimal temperature $\hat{T}^*$ with BN is superior to that without BN. These results validate the effectiveness of our proposed method.

## 5.2 Effect of the task

Revisiting Fig. 6, we can observe that while BN insertion can suppress the effect of differences in model architecture, the effect of differences in dataset slightly remains. We hypothesized that task difficulty causes this phenomenon. As the number of similar classes increases (i.e., the task difficulty increases), the greater the impact on learning since the

11

Table 2: Temperature determination coefficient estimated through optimization and the corresponding average accuracy [%].

| | w/o BN | | | w/ BN | | | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| | $\alpha$ | $\beta$ | $\mu_{acc}$ | $\alpha$ | $\beta$ | $\mu_{acc}$ | $\Delta_{acc}$ |
| **CIFAR10** | | | | | | | |
| VGG9 | 2.09 | -12.70 | 88.73 | 3.42 | -19.32 | 88.76 | 0.03 |
| ResNet10 | 7.62 | -60.92 | 90.74 | 0.84 | -4.29 | 91.11 | 0.37 |
| PyramidNet10 | 1.98 | -13.84 | 83.69 | 7.09 | -52.70 | 83.69 | 0.01 |
| **CIFAR100** | | | | | | | |
| VGG9 | 7.63 | -58.14 | 82.71 | 6.69 | -45.52 | 82.71 | 0.00 |
| ResNet10 | 3.79 | -30.07 | 85.96 | 2.68 | -17.42 | 86.41 | 0.45 |
| PyramidNet10 | 3.41 | -13.13 | 77.93 | 4.11 | -29.07 | 78.03 | 0.09 |
| **STL10** | | | | | | | |
| VGG9 | 3.10 | -22.80 | 79.78 | 2.95 | -15.30 | 80.41 | 0.63 |
| ResNet10 | 1.41 | -11.17 | 81.14 | 3.18 | -21.39 | 82.62 | 1.48 |
| PyramidNet10 | 4.64 | -20.47 | 74.08 | 2.98 | -7.87 | 74.46 | 0.38 |
| Avg. | 3.96 | -27.03 | 82.75 | 3.77 | -23.65 | 83.13 | 0.38 |
| S.D. | 2.17 | 18.29 | 4.95 | **1.86** | **15.30** | 4.89 | 0.44 |

output distribution is gradually flattened with increasing $T$, as illustrated in Fig. 1. Therefore, we analyzed changes in the temperature determination coefficients with respect to task difficulty. We constructed datasets with different task difficulties by extracting subsets containing 10 classes from CIFAR100 (referred to as the CIFAR100 subset). From these subsets, we utilized two relatively difficult subsets and two relatively easy subsets. Moreover, we simultaneously utilized the original CIFAR-100 and Tiny IN datasets. To calculate task difficulty, we employed the CSG [11]. We computed the CSG scale using the original images. While the original implementation trains an autoencoder and employs t-SNE, we do not use either of them, but it is sufficient for this case. For CSG calculation, we employed the publicly available implementation, which can be found at `https://github.com/Dref360/spectral-metric`.

In Fig. 9, we present the results on four subsets of CIFAR100, the original CIFAR100, and Tiny IN, employing ResNet10 as the model architecture. We observe the following trends, which are consistent with the previous section:

1. BN insertion stabilizes the temperature determination coefficients.

2. While BN insertion reduces the test accuracy at $T = 1$, it improves the test accuracy at $T = \hat{T}^*$.

3. Higher task difficulty scores (CSG) correlate with greater improvement due to temperature adjusting.

Furthermore, we observe that as the CSG increases, the optimal temperature $T^*$ decreases. An increase in the CN may be responsible for the decrease in $T^*$. These findings are consistent with the hypothesis that increased $T$ has a larger impact on similar classes, particularly in challenging tasks.

Based on the above, the optimal temperature determination coefficients can be estimated by inserting BN, representing $\hat{T}^*$ as a function of $M$, and correcting with CSG and CN. Based on Eq. (5), we defined as follows:

$$\hat{T}^* = \text{clip}(\alpha\sqrt{M} + \beta + \gamma\log(csg) + \delta\log(cn), \epsilon, 512), \tag{16}$$

In the case without correction, $\gamma = \delta = 0$. When only the CSG is used for correction (CSG Corr.), $\delta = 0$, while when only the CN is used for correction (CN Corr.), $\gamma = 0$. The $\text{clip}()$ function represents a clipping operation that limits the values between $\epsilon$ and 512. We experimentally set the upper bound of $T$ to 512 and the lower bound to 1.

As in the previous section, we performed optimization by cross-validation (CO) to estimate the temperature determination coefficients. Table 3 displays the test accuracy obtained by performing linear interpolation on each dataset. We observe a small improvement in overall performance, particularly notable on the Tiny IN dataset. For the original CIFAR100 dataset, while the performance decreases slightly, it can be improved by increasing the number of samples for simulation and optimization.

Table 3 also displays the results of performing global optimization (GO) on all experimental results displayed in Fig. 9. As expected, the average performance exceeds that obtained by optimization using cross-validation. While performing
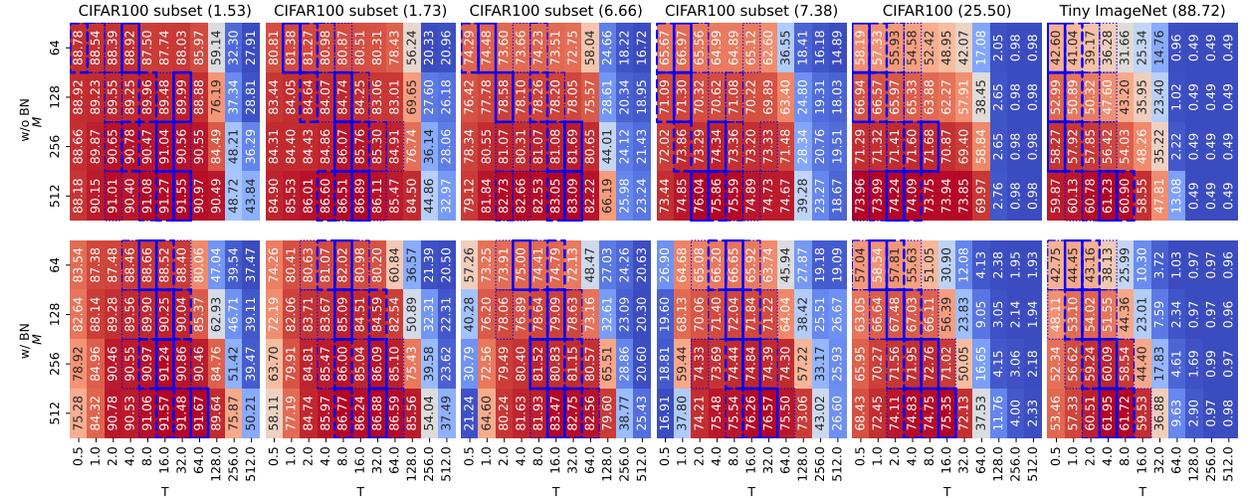
Figure 9: Test accuracies [%] for each temperature parameter in different task difficulty scenarios using ResNet10. The upper row displays results without BN insertion (conventional), while the lower row displays results with BN insertion (proposed method). Values in brackets indicate CSGs.

Table 3: Average estimation accuracy [%] achieved by our proposed methods. Each value indicates the average test accuracy for all $M$.

| Dataset | CIFAR100 subset | | | | CIFAR100 | Tiny IN | |
|---|---|---|---|---|---|---|---|
| CSG | 1.53 | 1.73 | 6.66 | 7.38 | 25.50 | 88.72 | Avg. |
| CN | 10 | 10 | 10 | 10 | 100 | 200 | |
| ResNet 9 w/ BN (CO) | 89.75 | 84.12 | 78.41 | 71.15 | 68.16 | 49.55 | 73.52 |
| + CSG Corr. (CO) | 90.26 | 84.66 | 79.18 | 72.05 | 66.43 | 54.48 | 74.51 |
| + CN Corr. (CO) | 90.31 | 84.51 | 79.29 | 72.29 | 67.56 | 54.30 | 74.71 |
| + CSG&CN Corr. (CO) | 90.20 | 84.61 | 79.48 | 72.33 | 66.58 | 54.66 | 74.64 |
| + CSG Corr. (GO) | 90.32 | 84.65 | 79.32 | 72.20 | 67.58 | 54.90 | 74.83 |
| + CN Corr. (GO) | 90.24 | 84.76 | 79.38 | 72.29 | 67.87 | 55.01 | 74.92 |
| + CSG&CN Corr. (GO) | 90.36 | 84.74 | 79.39 | 72.23 | 68.14 | 55.03 | 74.98 |

correction is beneficial, the difference in performance between CSG Corr. and CN Corr. is negligible. Therefore, one can simply use CN Corr. or CSG & CN Corr. when the CSG value can be computed. Finally, we can estimate the optimal temperature parameter $\hat{T}^*$ based on the temperature determination coefficients estimated by global optimization, as demonstrated in Table 4.

Table 4: Temperature determination coefficients obtained by global optimization.

| | $\alpha$ | $\beta$ | $\gamma$ | $\delta$ |
|---|---|---|---|---|
| $\hat{T}^*$ | 0.7239 | -4.706 | | |
| $\hat{T}^*_{csg}$ | 0.4111 | 6.848 | -2.024 | |
| $\hat{T}^*_{cn}$ | 0.4051 | 6.656 | | -1.973 |
| $\hat{T}^*_{csgcn}$ | 0.3192 | 20.74 | 3.746 | -7.38 |

13

# 6 Effectiveness validation with real problems

In the previous section, we discussed the performance of temperature parameters based on estimated test accuracy through linear interpolation. We also employed relatively shallow model architectures. To validate the robustness of our proposed method, we estimate the optimal temperature $\hat{T}^*$ for standard deep learning models. Based on the empirical data obtained during actual training and validation, we evaluate the effectiveness of our proposed method. The training settings are identical to those described in the previous section, with the exception of batch size, which varies across datasets.

We compared the following approaches:

1. $T = 1$: The conventional method that uses $T = 1$ without BN insertion (current **default** setting).

2. $T = \sqrt{M}$: The method used in self-attention [21], which employs $T = \sqrt{M}$ without BN insertion.

3. **Insert LN**: The method that inserts LN just before applying the softmax function without BN insertion [9].

4. **Ours** ($\hat{T}^*$): The method that inserts BN just before the output layer and employs $\hat{T}^*$.

5. **Ours** ($\hat{T}^*_{csg}$): The method that corrects the temperature parameter to $\hat{T}^*_{csg}$ calculated from the CSG with global optimization.

6. **Ours** ($\hat{T}^*_{cn}$): The method that correct the temperature parameter to $\hat{T}^*_{cn}$ based on the CN with global optimization.

7. **Ours** ($\hat{T}^*_{csgcn}$): The method that corrects the temperature parameter to $\hat{T}^*_{csgcn}$ using the CSG and CN with global optimization.

## 6.1 Evaluation of robustness to differences in model architecture

We validated our method on various model architectures, including ConvNets and ViTs.

### 6.1.1 Convolutional neural networks

We employed the following widely used ConvNets architectures: ResNet50 [32], EfficientNet (Eff. Net) b0/b5 [36], RegNet8GF [37], and ConvNeXt-Small [38]. We trained all models from scratch on CIFAR10 and CIFAR100.

In Table 5, we present the performance of each model. Notably, our proposed method outperforms other methods in many cases. Interestingly, $T = \sqrt{M}$ derived from $1/\sqrt{d_k}$ from the self-attention study [21] sometimes outperforms our method. This method differs from our proposed method in terms of BN insertion and the calculation of $T$. Since $T = \sqrt{M}$ does not significantly differ from our proposed method in some cases, it is expected to exhibit high performance, particularly for datasets with a low CSG. However, it is sensitive to changes in dataset and model architecture, resulting in significantly lower performance in cases such as ConvNeXt on CIFAR10 and Eff. Net b5 on CIFAR100. In terms of average performance, our method achieves the best performance on CIFAR10, while the addition of CSG Corr. achieves the best performance on CIFAR100. It should be noted that because CSG correction aims to improve robustness against difficult datasets such as Tiny ImageNet, it may not be as effective for easier datasets such as CIFAR10/100.

### 6.1.2 Vision transformers

We also evaluated our methods on recently developed ViT-based architectures: ViT-B/16 [39] and Swin Transformer Tiny (SwinT) [40]. We utilized Caltech101 [41] as the dataset.

In Table 6, we present the results for both models on the Caltech101 dataset. For both settings, our proposed method achieves the best performance. Although the performance significantly decreases with CSG and CN correction, it may be improved through individual hyperparameter tuning. For comparative experiments in a unified environment, we present the results as obtained.

## 6.2 Evaluation of robustness to dataset differences

We evaluated whether our proposed method was effective for various image classification datasets that were not used for optimization. We used the following popular vision datasets including fine-grained image recognition: German Traffic Sign Recognition Benchmark (GTS)[42], EuroSAT (SAT)[43], ISIC Challenge 2019 (ISIC)[44, 45, 46], Caltech101

Table 5: Performance evaluation for each Torchvision model in actual measurements [%].

| | ResNet50 | Eff. Net b0 | Eff. Net b5 | RegNet 8 GF | ConvNeXt-S | Avg. |
|---|---|---|---|---|---|---|
| **CIFAR10** | | | | | | |
| $T = 1$ | 84.80 | 87.10 | 88.83 | 88.72 | 70.02 | 83.89 |
| $T = \sqrt{M}$ | **88.70** | 87.51 | **91.01** | **91.18** | 68.07 | 85.29 |
| Insert LN [9] | 85.66 | 87.45 | 81.27 | 88.97 | 70.03 | 82.68 |
| Ours ($\hat{T}^*$) | 88.57 | _87.61_ | _90.79_ | _90.74_ | **77.94** | **87.13** |
| Ours ($\hat{T}^*_{csg}$) | _88.62_ | 87.41 | 90.61 | 90.47 | _76.68_ | _86.76_ |
| Ours ($\hat{T}^*_{cn}$) | 88.45 | 87.60 | 90.15 | 90.23 | 71.87 | 85.66 |
| Ours ($\hat{T}^*_{csgcn}$) | 88.46 | **87.72** | 89.75 | 90.30 | 71.07 | 85.46 |
| **CIFAR100** | | | | | | |
| $T = 1$ | 49.18 | 58.75 | 62.34 | 56.52 | **50.34** | 55.43 |
| $T = \sqrt{M}$ | **60.99** | 46.23 | 15.97 | **65.16** | 32.62 | 44.19 |
| Insert LN [9] | 49.64 | 58.48 | 55.50 | 55.41 | _49.58_ | 53.72 |
| Ours ($\hat{T}^*$) | _59.95_ | 59.01 | 66.05 | _64.94_ | 35.25 | 57.04 |
| Ours ($\hat{T}^*_{csg}$) | 58.70 | _59.20_ | **68.20** | 63.60 | 41.70 | **58.28** |
| Ours ($\hat{T}^*_{cn}$) | 57.93 | 58.94 | _67.81_ | 63.48 | 41.00 | 57.83 |
| Ours ($\hat{T}^*_{csgcn}$) | 57.89 | **60.14** | 65.93 | 63.56 | 42.70 | _58.04_ |

(CT)[41], Flowers102 (FLW)[47], Oxford-IIIT Pet (PET)[48], Describable Textures Dataset (DTD)[49], CUB-2011-200 dataset (CUB)[50]. The validation was conducted under the following three scenarios: (1) training ResNet50 from scratch (init), (2) fine-tuning a ResNet50 pre-trained on ImageNet-1k (FT), and (3) fine-tuning a SwinT model pre-trained on ImageNet-1k (FT). During the fine-tuning phase, both models were trained for 50 epochs using a maximum learning rate of 0.01.

Table 6: Performance evaluation of transformer models on the Caltech101 dataset in actual measurements [%].

| | ViT B/16 | SwinT |
|---|---|---|
| $T = 1$ | 50.38 | 54.07 |
| $T = \sqrt{M}$ | 47.54 | 61.29 |
| Insert LN [9] | 49.92 | 68.20 |
| Ours ($\hat{T}^*$) | 54.15 | 70.66 |
| Ours ($\hat{T}^*_{csg}$) | _54.30_ | _72.77_ |
| Ours ($\hat{T}^*_{cn}$) | **55.34** | **74.46** |
| Ours ($\hat{T}^*_{csgcn}$) | 5.03 [2] | 40.51 [3] |

In Table 7, we present the results for each dataset. Despite each correction formula $\hat{T}^*$ to $\hat{T}^*_{csgcn}$ being derived from CIFAR10/100 and Tiny ImageNet, our proposed method exhibits superior performance. Specifically, corrections based on CN or CSG&CN are effective in numerous instances, illustrating the broad applicability of our method to common vision tasks—even when those tasks were not used in the optimization of temperature determination coefficients.

Across the three scenarios, no significant differences in the observed trends were detected, indicating that the proposed method remains effective regardless of changes in model architecture or the approach used for training (i.e., from initial parameters or through fine-tuning). We initially hypothesized that the proposed method might be less effective if the model had already converged near a local optimum. However, even when applying the temperature determination coefficients optimized under scratch scenarios, the results confirmed that the proposed method remained effective for transfer-learning scenarios. These findings indicate that our proposed method can be effectively combined with transfer learning, preserving its benefits even in fine-tuning scenarios.

Examining the results across all datasets indicates that the proposed method generally achieves superior performance. However, it appears to be less effective, specifically on the ISIC dataset. To explore possible reasons for this, we

---

[2]53.61% when the maximum learning rate is changed to 0.01

[3]73.52% when the maximum learning rate is changed to 0.01

Table 7: Performance evaluation of three scenarios for each dataset in actual measurements [%].

| Dataset | GTS | SAT | ISIC | CT | FLW | PET | DTD | CUB | Avg. |
|---|---|---|---|---|---|---|---|---|---|
| CSG | 1.39 | 3.08 | 3.85 | 5.59 | 7.47 | 22.80 | 25.86 | 75.32 | |
| CN | 43 | 10 | 8 | 101 | 102 | 37 | 47 | 200 | |
| **ResNet50 init** | | | | | | | | | |
| $T = 1$ | 70.16 | 97.33 | **76.16** | 83.49 | 39.93 | 55.74 | 34.10 | 51.98 | 63.61 |
| $T = \sqrt{M}$ | 88.69 | 97.56 | 74.07 | 86.37 | 39.47 | 76.02 | 34.63 | 35.42 | 66.53 |
| Insert LN [9] | 71.56 | 96.85 | 75.17 | 87.63 | 48.84 | 71.46 | 36.65 | 54.92 | 67.89 |
| Ours ($\hat{T}^*$) | 87.26 | 97.59 | 66.57 | 87.52 | 53.54 | 76.12 | **42.29** | 63.01 | 71.74 |
| Ours ($\hat{T}^*_{csg}$) | 87.56 | 97.69 | 72.16 | 88.63 | **54.81** | **77.49** | 36.76 | 60.08 | 71.90 |
| Ours ($\hat{T}^*_{cn}$) | 92.01 | **97.90** | 70.03 | **90.59** | 54.76 | 74.03 | 35.53 | **64.29** | **72.39** |
| Ours ($\hat{T}^*_{csgcn}$) | **92.18** | 97.40 | 69.32 | 90.05 | 52.98 | 75.63 | 36.28 | 63.34 | 72.15 |
| **ResNet50 FT** | | | | | | | | | |
| $T = 1$ | 89.14 | 98.44 | 84.72 | **96.43** | 90.18 | 87.44 | 62.98 | 75.03 | 85.54 |
| $T = \sqrt{M}$ | 87.96 | 98.52 | 82.75 | 33.37 | 17.30 | 77.19 | 34.79 | 5.14 | 54.63 |
| Insert LN [9] | 91.84 | **98.58** | **84.93** | 95.31 | 90.32 | 88.09 | 63.56 | 75.13 | 85.97 |
| Ours ($\hat{T}^*$) | 92.43 | 98.46 | 84.39 | 93.93 | 79.10 | 89.59 | 64.20 | 66.55 | 83.58 |
| Ours ($\hat{T}^*_{csg}$) | 92.39 | 98.51 | 80.58 | 95.62 | 84.45 | 88.83 | 65.05 | **76.92** | 85.29 |
| Ours ($\hat{T}^*_{cn}$) | 92.29 | 98.43 | 83.88 | 96.08 | 89.17 | 89.26 | **65.64** | 75.46 | 86.28 |
| Ours ($\hat{T}^*_{csgcn}$) | **92.94** | 98.47 | 84.38 | 96.20 | **91.27** | **89.81** | 65.21 | 76.49 | **86.85** |
| **SwinT FT** | | | | | | | | | |
| $T = 1$ | 89.64 | **98.57** | 84.36 | 95.85 | 80.29 | 89.40 | 65.11 | 76.77 | 85.00 |
| $T = \sqrt{M}$ | 92.78 | 98.41 | 83.46 | 95.01 | 51.02 | 91.39 | 65.37 | 67.03 | 80.56 |
| Insert LN [9] | 91.03 | 98.25 | **84.51** | 96.85 | 77.18 | 91.20 | **68.51** | 74.53 | 85.26 |
| Ours ($\hat{T}^*$) | 92.49 | 98.41 | 80.88 | 93.74 | 66.47 | 91.71 | 66.54 | 57.40 | 80.96 |
| Ours ($\hat{T}^*_{csg}$) | 92.60 | 98.43 | 83.64 | 93.93 | 60.81 | 91.14 | 68.14 | 77.99 | 83.34 |
| Ours ($\hat{T}^*_{cn}$) | **93.24** | 98.43 | 80.83 | 96.66 | 76.57 | **92.15** | 67.61 | 79.19 | 85.58 |
| Ours ($\hat{T}^*_{csgcn}$) | 92.13 | 98.43 | 81.32 | **96.93** | **90.21** | 91.63 | 68.14 | **79.67** | **87.31** |

conducted performance evaluations on the ISIC dataset by varying the temperature and comparing conditions with and without BN in the ResNet50 transfer scenario, as shown in Fig. 10. Under the with-BN condition, the proposed method achieves results close to the optimal value, $\hat{T}_{csgcn} = 24.88$, confirming its efficacy. In contrast, under the without-BN condition, a wider temperature range yields high performance, resulting in comparatively poorer outcomes for the proposed method. As described in Section 4, incorporating BN generally improves performance when the optimal temperature is selected, although this pattern was not observed with the ISIC dataset. Further investigation is needed to clarify the conditions under which this discrepancy arises.

## 6.3   Performance evaluation with the incorporation of label smoothing

In Section 2.3, we discuss the difference between label smoothing and our proposed method. However, the effect of combining label smoothing remains unclear. Therefore, we evaluated the scratch and transfer performance of ResNet50 using label smoothing.

We present the results using label smoothing in Table 8. It can be seen that our proposed method performs effectively without label smoothing. However, when incorporating label smoothing, the performance of our proposed method decreases. In comparison with Table 7, the conventional method ($T = 1$) exhibits improved performance by incorporating label smoothing. From the discussion in Section 2.3, we expect both label smoothing and temperature scaling to play distinct roles in regularization. However, the experimental results indicate that this study's estimation formula for $T^*$ is incompatible with label smoothing. Nevertheless, since the proposed method outperforms label smoothing when $T = 1$ in terms of average performance, we conclude that it is preferable to use the proposed method rather than label smoothing.
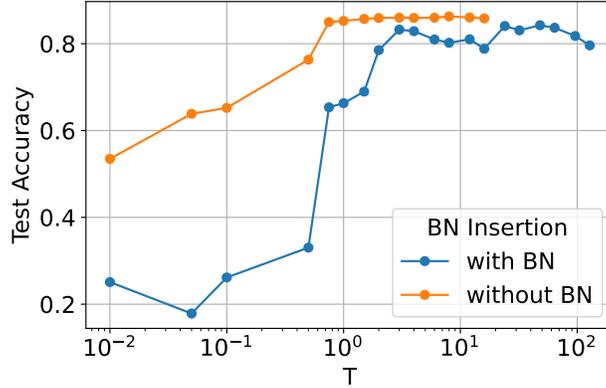
Figure 10: Test accuracies [%] for each temperature parameter in ISIC using ResNet50.

Table 8: Performance evaluation using the proposed method with label smoothing in ResNet50 [%].

| | GTS | CT | FLW | PET | DTD | Avg. |
|---|---|---|---|---|---|---|
| **Scratch** | | | | | | |
| $T = 1$ | 90.20 | 89.06 | 45.03 | 59.47 | 30.16 | 62.78 |
| $T = \sqrt{M}$ | 91.02 | 87.02 | 39.71 | **77.43** | **40.74** | 67.19 |
| Insert LN [9] | 91.64 | 88.75 | 50.59 | 72.64 | <u>39.57</u> | 68.64 |
| Ours ($\hat{T}^*$) | <u>91.97</u> | 90.02 | **57.42** | <u>76.42</u> | 36.70 | **70.51** |
| Ours ($\hat{T}^*_{csg}$) | 91.27 | 90.02 | <u>55.68</u> | 75.28 | 37.82 | 70.01 |
| Ours ($\hat{T}^*_{cn}$) | **92.00** | 89.78 | 54.89 | 75.12 | 39.04 | <u>70.17</u> |
| Ours ($\hat{T}^*_{csgcn}$) | 91.76 | **90.59** | 50.95 | 74.52 | 35.59 | 68.68 |
| **Transfer** | | | | | | |
| $T = 1$ | 90.17 | <u>96.31</u> | 90.55 | <u>90.19</u> | **65.85** | 86.61 |
| $T = \sqrt{M}$ | 86.04 | 34.79 | 16.44 | 62.47 | 33.67 | 46.68 |
| Insert LN [9] | <u>92.13</u> | 96.16 | <u>90.76</u> | 89.51 | 64.52 | <u>86.62</u> |
| Ours ($\hat{T}^*$) | 92.03 | 93.01 | 76.73 | 89.34 | 65.27 | 83.28 |
| Ours ($\hat{T}^*_{csg}$) | 91.39 | 95.47 | 83.57 | 90.11 | <u>65.59</u> | 85.22 |
| Ours ($\hat{T}^*_{cn}$) | **92.71** | 95.81 | 88.08 | 90.05 | 64.95 | 86.32 |
| Ours ($\hat{T}^*_{csgcn}$) | 91.69 | **96.39** | **91.33** | <u>90.24</u> | 65.05 | **86.94** |

## 7   Conclusion

In this study, we investigated the optimization of the temperature parameter in the commonly used softmax cross-entropy loss for classification problems using deep learning. It has been suggested that there is an important relationship between the number of dimensions $M$ of the encoder's feature map and the temperature $T$, and we aimed to theoretically verify this relationship and experimentally determine the optimal temperature $T^*$. Theoretical verification led to the hypothesis that standardizing the feature maps using BN enables robust calculation of the optimal temperature $T^*$ across different models and datasets. The experimental results support this hypothesis, demonstrating that BN insertion can lead to optimal performance, whereas using $T = 1$ often results in performance degradation.

Based on a comprehensive evaluation, we propose a method for estimating $T^*$ through optimization, along with a correction approach based on task difficulty (denoted by CSG) and the number of classes (denoted by CN), thereby yielding an optimal temperature $\hat{T}^*$ that leads to performance improvement. Ultimately, we derived the formula

$$\hat{T}^*_{csgcn} = \text{clip}\big(0.3192\sqrt{M} + 20.74 + 3.746\log(csg) - 7.380\log(cn), \epsilon, 512\big).$$

Evaluations across various model architectures and image classification datasets demonstrated that the proposed method consistently achieves high accuracy in most cases.

This study has several limitations. First, it focuses solely on image classification datasets. The related work [9] has demonstrated that similar phenomena occur in one-dimensional waveform activity recognition. This suggests that the proposed method may also be effective in other tasks employing the softmax cross-entropy loss, such as speech and document classification. Second, there is room for improvement in the correction formula from ($\hat{T}^*$ to $\hat{T}^*_{csgcn}$). The temperature determination coefficients were optimized based on CIFAR and Tiny ImageNet experiments. Therefore, diversifying measurement points and environments may lead to more generalized temperature determination coefficients. Currently, the temperature determination coefficients are optimized using data from up to 200 classes in the Tiny ImageNet dataset. Consequently, their effectiveness cannot be assured for scenarios involving a larger number of classes or substantially larger datasets. Furthermore, the correction for task difficulty was performed using the CSG, which involves some ambiguity and uncertainty. Therefore, we plan to further explore dynamic adjustments to task difficulty based on interclass similarity during training. Third, the temperature $T$ was kept constant in this study. Given the characteristics of $T$, it may be preferable for $T$ to be lower during the early stages of training and increase as the training progresses. We plan to investigate the implications of a dynamically changing $T$ during training in the future work.

## Appendix A    Inserting BN

Normalization layers, such as BN [12] and LN [20], aim to normalize feature maps. These layers mitigate the distribution shift between layers, referred to as the internal covariate shift. These normalization layers have been introduced by recent architectures, such as ViT [39] and the Swin transformer (Swin-T) [40]. However, it is uncommon to insert normalization directly before the output layer, as proposed in this study. Referring to official PyTorch implementations[4], many models, such as ResNet [32], RegNet [37], and EfficientNet [36], typically connect the output layer after an encoder. Among the models we examined, ConvNeXt [38], ViT and Swin-T connect the output layer following the sequence: encoder $\rightarrow$ GAP $\rightarrow$ LN.

Based on our hypothesis described in the previous section 3.2.2, we employ BN as the normalization layer, as it normalizes features for the number of samples. Inserting a normalization layer immediately before the output layer has not traditionally been a common practice and has only been adopted in certain specific model architectures mentioned above. While the insertion of BN is pointed out to have a stabilizing effect on temperature determination coefficients, there are concerns that it may degrade the original performance, necessitating verification (we discuss in Section 4).

## Appendix B    Additional Theoretical Justification of the Temperature Coefficients

This appendix details the derivation of the two task–dependent correction terms—CSG correction and class-number (CN) correction—that appear in Eq. (5).

### B.1    CSG Correction Term

With BN inserted immediately before the output layer, the logit variance is given by Eq. (15), reproduced here for convenience:

$$\mathbb{V}[\hat{y}_j] = \sum_{k=1}^{M} w_{jk}^2 + \sum_{k \neq l} w_{jk} w_{jl} \, \mathbb{E}[z_k z_l]. \tag{A.2}$$

The first summation depends only on the architecture, whereas the second captures the residual inter-feature correlations. Empirically, these correlations increase with task difficulty, and can be approximated by

$$\sum_{k \neq l} w_{jk} w_{jl} \, \mathbb{E}[z_k z_l] \; \propto \; \log(\text{CSG}). \tag{A.3}$$

Because a larger logit variance flattens the soft-max distribution—effectively acting as a higher temperature—the temperature must be lowered in accordance with task difficulty to counter this effect. We therefore adopt CSG as a proxy measure of task difficulty.

### B.2    Class-Number (CN) Correction Term

Let $\mathbf{y} = (y_1, \ldots, y_C)$ be row–wise i.i.d. logits and $p_i = e^{y_i/T} / \sum_{k=1}^{C} e^{y_k/T}$ the corresponding softmax probabilities. Because $\sum_{i=1}^{C} p_i = 1$ and the distribution of $\mathbf{y}$ is exchangeable, the unconditional expectation of any single class

---

[4]PyTorch: Models and pre-trained weights [`https://pytorch.org/vision/main/models.html`]

probability is $\mathbb{E}[p_i] = \frac{1}{C}$. This equation shows that the baseline confidence allocated to each class decays inversely with the number of classes $C$. Consequently, the maximum softmax probability in a sample $\max_j p_j$ also decreases monotonically with $C$, a tendency confirmed by the corrected simulation in Fig. 11.
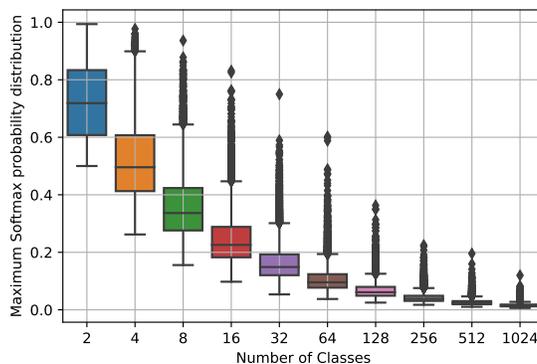


Figure 11: Distribution of the maximum softmax probability for varying class counts.

For the cross-entropy loss $L = -\log p_{y^*}$ ($y^*$ is the ground-truth label), the gradient with respect to a logit $z_i$ is

$$\frac{\partial L}{\partial z_i} = \frac{1}{T}\big(p_i - \mathbb{1}_{\{i=y^*\}}\big).$$

If $C$ grows while the logit scale remains unchanged, $p_{y^*}$ converges to $1/C$ and $\partial L/\partial z_i$ shrinks approximately as $1/T\,C$. The optimisation thus suffers from vanishing gradients. One remedy is to reduce $T$ (i.e. sharpen the softmax) so that the effective gradient magnitude stays within a favourable range. Hence, $T^*$ must decrease in expectation with $C$.

Because the effect described above is systematic and orthogonal to the feature-dimension effect (first term), we embed $C$ in the model through the logarithmic correction term $\delta \log(\mathrm{cn})$ in Eq. (5). In Fig. 11, we show the distribution of the maximum softmax probability $\max_j p_j$ as a function of the number of classes $C$. For clarity, the horizontal axis is plotted on a $\log_2$ scale so that the powers of two $(2, 4, 8, \cdots, 1024)$ appear evenly spaced. This choice is purely empirical and serves to improve readability; it is not derived from any grid-search procedure. The plot confirms that $\max_j p_j$ decays monotonically toward zero as $C$ increases, although the decay is not strictly linear in $\log C$. Optimising the coefficient $\delta$ across diverse datasets yields the globally valid estimate reported in Table 4.

## Acknowledgement

**Availability of data and material** We used publicly available datasets CIFAR10/100 [28], STL10 [29], Tiny ImageNet (Tiny IN) [30], German Traffic Sign Recognition Benchmark (GTS)[42], EuroSAT (SAT)[43], ISIC Challenge 2019 (ISIC)[44, 45, 46], Caltech101 (CT)[41], Flowers102 (FLW)[47], Oxford-IIIT Pet (PET)[48], Describable Textures Dataset (DTD)[49], CUB-2011-200 dataset (CUB)[50] for our experiments. All datasets are available for each web site.

**Competing interests** The authors declare that they have no conflicts of interest.

## References

[1] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 2012.

[2] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg. Ssd: Single shot multibox detectors. *European conference on computer vision*, pages 21–37, 2016.

[3] T. Karras, S. Laine, and T. Aila. A stylebased generator architecture for generative adversarial networks. *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019.

[4] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Nee-lakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901, 2020.

[5] G. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. *NIPS Deep Learning and Representation Learning Workshop (online)*, 2015.

[6] C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger. On calibration of modern neural networks. *Proceedings of the 34th International Conference on Machine Learning*, 70:1321–1330, 2017.

[7] Hassan Ismail Fawaz, Germain Forestier, Jonathan Weber, Lhassane Idoumghar, and Pierre-Alain Muller. Deep learning for time series classification: a review. *Data Min. Knowl. Discov.*, 33(4):917–963, July 2019.

[8] Atish Agarwala, Samuel Stern Schoenholz, Jeffrey Pennington, and Yann Dauphin. Temperature check: theory and practice for training models with softmax-cross-entropy losses. *Transactions on Machine Learning Research*, 2023.

[9] Tatsuhito Hasegawa. Optimal temperature parameter of softmax while training deep learning model in activity recognition. *Journal of Information Processing Society of Japan*, 64(8):1182–1192, August 2023. (in Japanese).

[10] Hao Xuan, Bokai Yang, and Xingyu Li. Exploring the impact of temperature scaling in softmax for classification and adversarial robustness, 2025.

[11] Frédéric Branchaud-Charron, Andrew Achkar, and Pierre-Marc Jodoin. Spectral metric for dataset complexity assessment. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3210–3219, 2019.

[12] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *ICML*, 2015.

[13] Jihao Liu, Boxiao Liu, Hongsheng Li, and Yu Liu. Meta knowledge distillation, 2022.

[14] Zheng Li, Xiang Li, Lingfeng Yang, Borui Zhao, Renjie Song, Lei Luo, Jun Li, and Jian Yang. Curriculum temperature for knowledge distillation. In *Proceedings of the Thirty-Seventh AAAI Conference on Artificial Intelligence*, 2023.

[15] Shangquan Sun, Wenqi Ren, Jingzhi Li, Rui Wang, and Xiaochun Cao. Logit standardization in knowledge distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15731–15740, 2024.

[16] Sergio A. Balanya, Juan Maroñas, and Daniel Ramos. Adaptive temperature scaling for robust calibration of deep neural networks. *Neural Computing and Applications*, 36(14):8073–8095, May 2024.

[17] Feng Wang and Huaping Liu. Understanding the behaviour of contrastive loss. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2495–2504, 2021.

[18] Lukás Neumann, Andrew Zisserman, and A Vedaldi. Relaxed softmax: Efficient confidence auto-calibration for safe pedestrian detection. In *Proceedings of the 2018 NIPS Workshop on Machine Learning for Intelligent Transportation Systems*, 2018.

[19] Xu Zhang, Felix Xinnan Yu, Svebor Karaman, Wei Zhang, and Shih-Fu Chang. Heated-up softmax embedding, 2018.

[20] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization, 2016. arXiv:1607.06450.

[21] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS'17, page 6000–6010, 2017.

[22] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna. Rethinking the inception architecture for computer vision. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2818–2826, 2016.

[23] Rafael Müller, Simon Kornblith, and Geoffrey E Hinton. When does label smoothing help? In *Advances in Neural Information Processing Systems*, volume 32, 2019.

[24] Chang-Bin Zhang, Peng-Tao Jiang, Qibin Hou, Yunchao Wei, Qi Han, Zhen Li, and Ming-Ming Cheng. Delving deep into label smoothing. *Trans. Img. Proc.*, 30:5984–5996, jan 2021.

[25] Yann LeCun, Léon Bottou, Genevieve B. Orr, and Klaus-Robert Müller. Efficient backprop. In *Neural Networks: Tricks of the Trade, This Book is an Outgrowth of a 1996 NIPS Workshop*, page 9–50, Berlin, Heidelberg, 1998. Springer-Verlag.

[26] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, volume 9, pages 249–256, 2010.

[27] K. He, X. Zhang, S. Ren, and J. Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2015.

[28] K. Alex and H. Geoffrey. Learning multiple layers of features from tiny images. Master's thesis, University of Toronto, 2009.

[29] Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised feature learning. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, volume 15, pages 215–223, 2011.

[30] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. *CVPR09*, 2009.

[31] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *Proceedings of the International Conference on Learning Representations*, pages 1–14, 2015.

[32] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.

[33] D. Han, J. Kim, and J.n Kim. Deep pyramidal residual networks. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR*, pages 6307–6315, 2017.

[34] I. Loshchilov and F Hutter. Sgdr: Stochastic gradient descent with warm restarts. *International Conference on Learning Representations*, 2017.

[35] Rainer Storn and Kenneth Price. Differential evolution – a simple and efficient heuristic for global optimization over continuous spaces. *Journal of Global Optimization*, 11(4):341–359, Dec 1997.

[36] Mingxing Tan and Quoc Le. EfficientNet: Rethinking model scaling for convolutional neural networks. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 6105–6114, 2019.

[37] Ilija Radosavovic, Raj Prateek Kosaraju, Ross Girshick, Kaiming He, and Piotr Dollar. Designing network design spaces. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.

[38] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11976–11986, June 2022.

[39] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021.

[40] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9992–10002, 2021.

[41] Li Fei-Fei, Rob Fergus, and Pietro Perona. Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories. *Computer vision and Image understanding*, 106(1):59–70, 2007.

[42] Sebastian Houben, Johannes Stallkamp, Jan Salmen, Marc Schlipsing, and Christian Igel. Detection of traffic signs in real-world images: The German Traffic Sign Detection Benchmark. In *International Joint Conference on Neural Networks*, number 1288, 2013.

[43] Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 2019.

[44] Philipp Tschandl, Cliff Rosendahl, and Harald Kittler. The ham10000 dataset, a large collection of multi-source dermatoscopic images of common pigmented skin lesions. *Scientific Data*, 5(1):180161, Aug 2018.

[45] Noel C. F. Codella, David Gutman, M. Emre Celebi, Brian Helba, Michael A. Marchetti, Stephen W. Dusza, Aadi Kalloo, Konstantinos Liopyris, Nabin Mishra, Harald Kittler, and Allan Halpern. Skin lesion analysis toward melanoma detection: A challenge at the 2017 international symposium on biomedical imaging (isbi), hosted by the international skin imaging collaboration (isic). In *2018 IEEE 15th International Symposium on Biomedical Imaging (ISBI 2018)*, pages 168–172, 2018.

[46] Carlos Hernández-Pérez, Marc Combalia, Sebastian Podlipnik, Noel C. F. Codella, Veronica Rotemberg, Allan C. Halpern, Ofer Reiter, Cristina Carrera, Alicia Barreiro, Brian Helba, Susana Puig, Veronica Vilaplana, and Josep Malvehy. Bcn20000: Dermoscopic lesions in the wild. *Scientific Data*, 11(1):641, Jun 2024.

[47] M-E. Nilsback and A. Zisserman. Automated flower classification over a large number of classes. In *Proceedings of the Indian Conference on Computer Vision, Graphics and Image Processing*, Dec 2008.

[48] Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and C. V. Jawahar. Cats and dogs. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 3498–3505, 2012.

[49] Mircea Cimpoi, Subhransu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing textures in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014.

[50] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The caltech-ucsd birds-200-2011 dataset. Technical Report CNS-TR-2011-001, California Institute of Technology, 2011.