

# Supplementary Material: Calibration Transfer via Knowledge Distillation

Ramya Hebbalaguppe<sup>1,2§</sup>, Mayank Baranwal<sup>2§</sup>, Kartik Anand<sup>1</sup>, Chetan Arora<sup>1</sup>

<sup>1</sup>Indian Institute of Technology, Delhi    <sup>2</sup>Tata Consultancy Services Research

\*\* source code: <https://github.com/rhebbalaguppe/CalibrationXferViaKD>

## 1 Introduction

We complement our main text with supplementary materials encompassing the following components:

1. **Theoretical Insights:** This section contains main theoretical results, along with an explanation of the rationale behind choosing Padé approximants over more commonly used Taylor approximation, as discussed in Lemma 1.
2. **Rationale for Enhanced Performance:** This section elucidates the superior performance of the KD(C) framework, attributing it to three key factors: (a) insights from penultimate visualizations, (b) considerations of inter-class semantic similarities, and (c) the careful design of calibrators for the teacher model.
3. **Illustration of Generality:** Included is Fig. S5, which provides a visual demonstration of KD(C)'s versatility by comparing direct calibration with the KD(C) framework. It also presents an example featuring the [5] regularizer.
4. **Expanded Experimental Scope:** We strengthen the KD(C) methodology with additional experiments, covering various scenarios, including large-to-small, small-to-large, self-distillation, and iterative self-distillation. These experiments involve different descriptors and datasets.
5. **Additional Results:** We provide supplementary results that encompass calibration performance in the presence of dataset drift and reliability diagrams featuring confidence histograms, as elaborated in Sec. 5.1.
6. **Hyperparameter Analysis:** A detailed study explores how calibration and accuracy are influenced by various hyperparameters in the KD(C) framework, as depicted in Fig. S8.
7. **Source Code:** The supplementary materials include the source code along with a `readme.md` file, enclosed within the provided zip file.
8. **Training and Compute Details:** We furnish comprehensive information on the specifics of training and compute resources employed in our experiments.
9. **Limitations and Broader Impact:** This section delves into the limitations of our research and contemplates its broader impact on the field.

These supplementary materials serve to enrich and provide a deeper understanding of our main findings and contributions.

---

\*\* Equal contribution

## 2 Theoretical support: Additional details

### 2.1 Proof of Theorem 1

The proof of Theorem 1 is contingent on several essential Lemmas, which will be introduced beforehand. Lemma 1 and Lemma 2 capture the effect of quadratic temperature scaling in the KD loss function,  $\mathcal{L}_{KD}$ . In particular, it is shown that the partial derivative of  $\mathcal{L}_{KD}$  w.r.t. student’s logit for a given sample is equal to the difference in predicted probabilities of the student and teacher classifiers for that sample. These results are leveraged to characterize the first-order condition of optimality for the total loss function  $\mathcal{L}_{\text{tot}}$  w.r.t. parameters of the student classifier.

**Lemma 1.** *Let  $z_{i,s} := \mathbf{W}_s^\top \mathbf{x}_i$  and  $z_{i,t} := \mathbf{W}_t^\top \mathbf{x}_i$  with  $\tilde{p}_{i,s}, \tilde{p}_{i,t}$  be defined as above. Then  $\lim_{T \rightarrow \infty} T(\tilde{p}_{i,s} - \tilde{p}_{i,t}) \approx p_{i,s} - p_{i,t}$ .*

*Proof.* The result follows as a consequence of Padé approximation. Recall that from definition,

$$\tilde{p}_{i,s} - \tilde{p}_{i,t} = \frac{1}{1 + e^{-z_{i,s}/T}} - \frac{1}{1 + e^{-z_{i,t}/T}} \approx \frac{1}{1 + \frac{1 - \frac{z_{i,s}}{2T}}{1 + \frac{z_{i,s}}{2T}}} - \frac{1}{1 + \frac{1 - \frac{z_{i,t}}{2T}}{1 + \frac{z_{i,t}}{2T}}}, \quad (\text{S1})$$

where the last approximation follows from Padé approximation of the exponential when  $T$  is large.

Thus, Eq. (S1) can be re-written as:

$$\tilde{p}_{i,s} - \tilde{p}_{i,t} = \frac{1 + z_{i,s}/2T}{2} - \frac{1 + z_{i,t}/2T}{2} \implies T(\tilde{p}_{i,s} - \tilde{p}_{i,t}) = \frac{z_{i,s} - z_{i,t}}{4}. \quad (\text{S2})$$

On the other hand, a similar analysis following the Padé approximation yields:

$$p_{i,s} - p_{i,t} \approx (z_{i,s} - z_{i,t})/4. \quad (\text{S3})$$

Thus, Lemma 1 follows directly from Eq. (S2) and Eq. (S3).

**Remark:** Padé approximants have a wider range of convergence than the corresponding Taylor series, and can even converge where the Taylor series does not. For a detailed exposition, please refer to Sec. 2.3 and Fig. S1.

The following result shows that the quadratic temperature scaling in the KD loss function ensures that the gradients used to update the network weights are independent of the smoothed labels.

**Lemma 2 (Quadratic temperature scaling).** *Let  $\mathcal{L}_{KD}$  be defined as in Eq. (2). Then,*

$$\lim_{T \rightarrow \infty} \frac{\partial \mathcal{L}_{KD}}{\partial z_{i,s}} = p_{i,s} - p_{i,t}.$$

*Proof.* Recall that by definition  $\tilde{p}_{i,s} = \frac{1}{1 + e^{-z_{i,s}/T}}$ . The partial derivative of  $\tilde{p}_{i,s}$  w.r.t.  $z_{i,s}$  reads:

$$\frac{\partial \tilde{p}_{i,s}}{\partial z_{i,s}} = \frac{1}{T} \tilde{p}_{i,s} (1 - \tilde{p}_{i,s}). \quad (\text{S4})$$

On the other hand,

$$\frac{\partial \mathcal{L}_{KD}}{\partial z_{i,s}} = -T^2 \left( \frac{\tilde{p}_{i,t}}{\tilde{p}_{i,s}} - \frac{1 - \tilde{p}_{i,t}}{1 - \tilde{p}_{i,s}} \right) \frac{\partial \tilde{p}_{i,s}}{\partial z_{i,s}} = T^2 \frac{(\tilde{p}_{i,s} - \tilde{p}_{i,t})}{\tilde{p}_{i,s}(1 - \tilde{p}_{i,s})} \frac{\partial \tilde{p}_{i,s}}{\partial z_{i,s}}. \quad (\text{S5})$$

Thus, from Eq. (S4), Lemma 1 and for large  $T$ , Eq. (S5) reduces to:

$$\lim_{T \rightarrow \infty} \frac{\partial \mathcal{L}_{KD}}{\partial z_{i,s}} = p_{i,s} - p_{i,t},$$

which completes the proof.

**Lemma 3.** The derivative of the total loss function  $\mathcal{L}_{tot}$  w.r.t. the parameters  $\mathbf{W}_s$  of the student network lies in the span of  $\mathbf{X}$ , and is given by:

$$\frac{\partial \mathcal{L}_{tot}}{\partial \mathbf{W}_s} = \sum_{i=1}^N (p_{i,s} - \{\alpha p_{i,t} + (1 - \alpha)y_i\}) \mathbf{x}_i.$$

*Proof.* The proof follows directly from Lemma 2.

**Theorem 1.** Let  $\mathbf{X} \in \mathbb{R}^{d \times N}$  be the data matrix, and  $\mathbf{W}_s$  and  $\mathbf{W}_t$  represent the parameters of the student and the teacher networks, respectively. Then, under Assumption 1 and using the gradient-descent algorithm, the parameters  $\mathbf{W}_s$  of the student network converge to:

$$\mathbf{W}_s \approx \begin{cases} \alpha \mathbf{W}_t + 4(1 - \alpha) \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{Y}_{1/2}, & \text{if } N < d \\ \alpha \mathbf{W}_t + 4(1 - \alpha) (\mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{X} \mathbf{Y}_{1/2}, & \text{else} \end{cases},$$

where  $\mathbf{Y}_{1/2} := [y_i - \frac{1}{2}]_{i=1}^N$  is an  $N$ -dimensional vector.

*Proof.* First observe that the minimum value of the total loss function in Eq. (2) is finite. Moreover, the total loss function is convex in the parameters of the student network. Thus, any gradient-based descent algorithm with suitable step-size will converge to the optimizer asymptotically fast.

We now characterize the set of optimizers. Recall that the first-order condition of optimality implies:

$$\frac{\partial \mathcal{L}_{tot}}{\partial \mathbf{W}_s} = 0 \implies \sum_{i=1}^N (p_{i,s} - \{\alpha p_{i,t} + (1 - \alpha)y_i\}) \mathbf{x}_i = 0,$$

where the last equality follows from Lemma 3. Since the vectors  $\{\mathbf{x}_i\}$  are linearly independent (see Remark 1), the above equality holds if:

$$p_{i,s} - \{\alpha p_{i,t} + (1 - \alpha)y_i\} = 0, \quad \forall i \in \{1, \dots, N\}. \quad (\text{S6})$$

Expanding Eq. (S6) in terms of logits  $z_{i,s}$  leads to:

$$\frac{1}{1 + e^{-z_{i,s}}} = \alpha \frac{1}{1 + e^{-z_{i,t}}} + (1 - \alpha)y_i \implies \frac{1}{1 + \frac{1 - \frac{z_{i,s}}{2}}{1 + \frac{z_{i,s}}{2}}} \approx \alpha \frac{1}{1 + \frac{1 - \frac{z_{i,t}}{2}}{1 + \frac{z_{i,t}}{2}}} + (1 - \alpha)y_i, \quad (\text{S7})$$

where the last equation follows from Padé approximation. Rearranging the terms in Eq. (S7), and using the fact that  $z_{i,s} = \mathbf{W}_s^\top \mathbf{x}_i$  and  $z_{i,t} = \mathbf{W}_t^\top \mathbf{x}_i$ , one obtains:

$$(\mathbf{W}_s - \alpha \mathbf{W}_t)^\top \mathbf{x}_i = 4(1 - \alpha)(y_i - 1/2).$$

Since the above condition holds for every  $i \in \{1, \dots, N\}$ , the vector form of it can be written as:

$$\mathbf{X}^\top (\mathbf{W}_s - \alpha \mathbf{W}_t) = 4(1 - \alpha) \mathbf{Y}_{1/2}, \quad (\text{S8})$$

which, for  $N < d$ , is an underdetermined system of linear equations whose least-norm solution is given by:

$$\mathbf{W}_s = \alpha \mathbf{W}_t + 4(1 - \alpha) \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{Y}_{1/2}. \quad (\text{S9})$$

On the other hand, for  $N > d$ , (S8) represents an overdetermined system of linear equations whose least-norm solution is given by:

$$\mathbf{W}_s = \alpha \mathbf{W}_t + 4(1 - \alpha) (\mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{X} \mathbf{Y}_{1/2}. \quad (\text{S9})$$

which completes the proof.

## 2.2 Proof of Theorem 2

**Theorem 2.** *Let Assumption 1 hold. Let  $t_c$  and  $t_{uc}$  be two teacher classifiers with output probabilities  $\{p_{i,t_c}\}$  and  $\{p_{i,t_{uc}}\}$ , respectively. Also, let  $s_c$ ,  $s_{uc}$  depict two student classifiers trained independently from the corresponding teacher classifiers  $t_c$  and  $t_{uc}$  through KD, with output probabilities  $\{p_{i,s_c}\}$  and  $\{p_{i,s_{uc}}\}$ , respectively. Furthermore, assume that the teacher classifier  $t_c$  is well calibrated, then the student classifier  $s_c$  is also well calibrated. Conversely, if the teacher classifier  $t_{uc}$  is uncalibrated, the corresponding student classifier  $s_{uc}$  mimics a similar behavior, i.e.,*

$$\sum_{i=1}^N p_{i,s_c} = \sum_{i=1}^N y_i, \quad \text{and} \quad \sum_{i=1}^N p_{i,s_{uc}} \neq \sum_{i=1}^N y_i.$$

*Proof.* From Eq. (S6), the first-order condition for optimality for a student  $s$  trained from a teacher  $t$  through KD reads:

$$\sum_{i=1}^N p_{i,s} = \alpha \sum_{i=1}^N p_{i,t} + (1 - \alpha) \sum_{i=1}^N y_i,$$

which can be rewritten as

$$\sum_{i=1}^N (p_{i,s} - y_i) = \alpha \sum_{i=1}^N (p_{i,t} - y_i).$$

Thus for the same value of  $\alpha \in (0, 1)$ , if the teacher classifier  $s_c$  is well calibrated, then

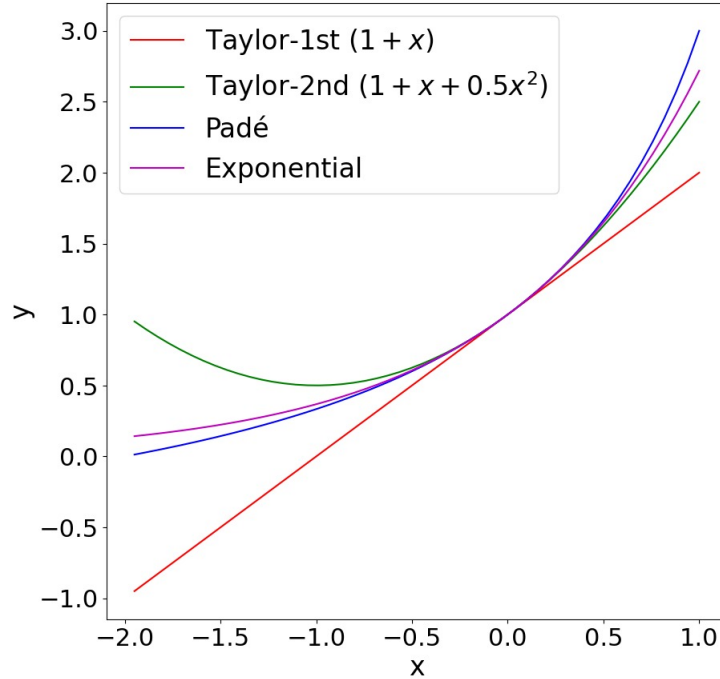
$$\sum_{i=1}^N (p_{i,s_c} - y_i) = \alpha \left( \sum_{i=1}^N (p_{i,t_c} - y_i) \right) = 0,$$

where, the last equality follows from well calibration of the teacher classifier. On the other hand,

$$\sum_{i=1}^N (p_{i,s_{uc}} - y_i) = \alpha \sum_{i=1}^N (p_{i,t_{uc}} - y_i) \neq 0 \implies \sum_{i=1}^N p_{i,s_{uc}} \neq \sum_{i=1}^N y_i,$$

which completes the proof.

### 2.3 Padé vs Taylor approximants



**Fig. S1:** Padé vs Taylor for a simple exponential function. Note that Padé approximants offer superior reliability compared to the extensively used Taylor approximants.

Employing approximants to derive theoretical outcomes in DNNs is commonplace due to the intricacies of dealing with highly nonlinear equations. We illustrate the difference between Padé and Taylor’s approximation as follows: Padé approximants have a wider range of convergence than the corresponding Taylor series, and can even converge where the Taylor series does not. A simple example of Padé approximant is,  $e^x = \frac{e^{0.5x}}{e^{-0.5x}} \approx \frac{(1+0.5x)}{(1-0.5x)} = (1+0.5x)(1-0.5x)^{-1}$ , which for  $|x| < 2$  can further be expanded to  $e^x \approx (1+0.5x)(1-0.5x)^{-1} = (1+0.5x)(1+0.5x+0.25x^2+\dots)$ . Thus, despite using first-order approximations for both the numerator and denominator terms, the above Padé approximant very closely follows the original exponential function. This is in contrast to Taylor’s expansion, and even a second-order Taylor’s expansion does not mimic the exponential function, except for a very small interval around the origin. Please refer to Figure S1 for further details.

This is precisely why we restrict using Padé approximants in our theoretical exploration since they are still potentially non-divergent in regimes even when  $z_{i,t}$  and  $z_{i,s}$  are not vanishingly small. It must also be remarked that **exact characterization of weights of student network is a theoretically hard problem**, and such practical approximations are useful to obtain important theoretical insights.

### 3 Rationale on the superior performance of our KD(C) framework

The essence of Theorem 2 lies in its assertion that uncalibrated teachers can only transfer their lack of calibration to their student counterparts, whereas calibrated teachers enable the distillation of calibrated students. This theorem underscores the crucial significance of utilizing calibrated teachers in the knowledge distillation process. In light of this observation, we advocate for a novel approach to achieving accurate and calibrated models: calibrating a model through distillation from another model that is already calibrated. To validate the efficacy of this approach, we conducted an extensive series of experiments, showcasing the capabilities of our framework, KD(C). Our experimental results provide compelling evidence that KD(C) yields student models characterized by two key attributes: dynamic calibration at the sample level and semantic calibration. These findings substantiate the effectiveness of our proposed framework in achieving both sample-level and semantic calibration in student models.

#### 3.1 Classification of label smoothing

**Standard/static label smoothing.** Label Smoothing (LS) serves as a regularization technique designed to address potential inaccuracies within datasets. It recognizes that maximizing the likelihood directly, denoted as  $P(y|\mathbf{x})$ , may be detrimental due to the possibility of errors in the training labels. To mitigate this issue, LS introduces controlled noise into the labeling process. In essence, LS operates as follows: Given a small constant value  $\epsilon$ , it considers the training label  $y$

to be correct with a probability of  $(1 - \epsilon)$  and incorrect otherwise. Specifically, in the context of a softmax model with  $k$  outputs, it replaces the traditional binary classification targets of 0 and 1 with modified targets. These modified targets consist of  $\frac{\epsilon}{(k-1)}$  for incorrect labels and  $(1 - \epsilon)$  for correct labels [12, 16]. This approach ensures that all output probabilities undergo uniform regularization, thereby helping to combat overfitting and improve model generalization.

**Adaptive Label Smoothing.** In this method the level of regularization applied to training labels, which are typically one-hot encoded, is dynamically adjusted based on the network’s output probabilities for different classes [2, 5, 13]. This method is found to be more beneficial than conventional static label smoothing (LS) proposed in [16].

**Conditional Label Smoothing.** In this method the training labels go through selective modifications based on specific criteria, such as the application of margin-based penalties [10]. This approach places its emphasis on and applies regularization solely to the probabilities that exhibit miscalibration, thereby demonstrating enhanced calibration capabilities.

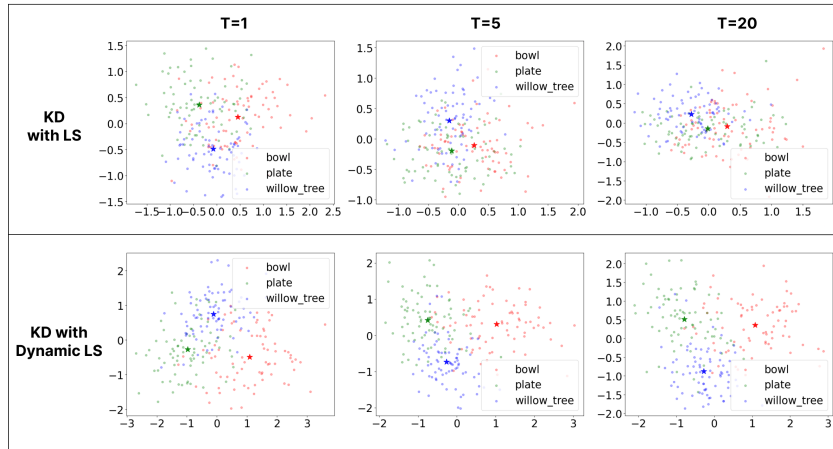
### 3.2 Visualization of penultimate layer’s activations reveal KD(C) using dynamic regularization works better than static regularization

**Penultimate Visualization.** [12] introduced this visualization technique wherein they projected the penultimate activations onto the hyperplane defined by the template vectors (weight vectors) corresponding to the selected classes (three classes) for visualization.

**Systematic diffusion.** The concept of “systematic diffusion”, introduced by [1], was developed to address discrepancies observed in prior studies, particularly the contradictions between [14] and the insights presented in LS literature [12]. This concept aims to elucidate the compatibility of label smoothing with knowledge distillation. The findings from [1]’s work indicate that when KD is conducted at elevated temperatures from a teacher model trained with LS, it results in a systematic shift in the relationships between classes. Specifically, for semantically similar classes, the inter-cluster distance decreases, while for the remaining classes, it increases relatively. Importantly, this diffusion of classes is not random; rather, it follows a systematic pattern.

In Fig. S2 and Fig. S3, we provide visual evidence of the limitations associated with LS-trained teachers compared to MDCA teachers [5]. These penultimate layer visualizations, inspired by the work of [14], reveal that semantically similar classes experience systematic diffusion when using LS, whereas this phenomenon is not observed with MDCA calibration. This observation substantiates our recommendation to opt for dynamic smoothing regularization techniques such as MDCA.

Notably, we notice a trend where distilled student models are most calibrated when the distillation temperature ( $T$ ) is approximately 1. We hypothesize that increasing  $T$  leads to the destruction of discriminating features, as outlined by [1],



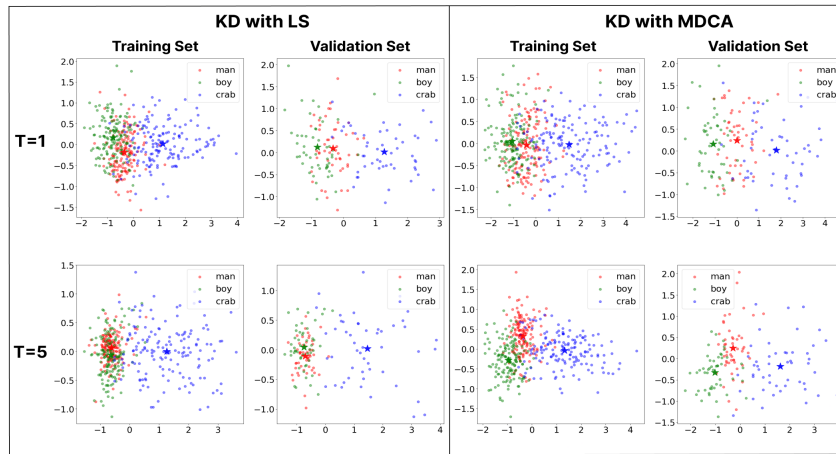
**Fig. S2: Visualization of penultimate layer’s activations** (Teacher = ResNet56, Student = ResNet8, Dataset = CIFAR100). We train ResNet8 using calibration techniques: KD with LS [16] (Left column) and KD with MDCA [5] Column). We follow the same setup and procedure used in papers [ [12, 14]] We use two semantically similar classes (**bowl**, **plate**) and one semantically dissimilar class (**willow\_tree**). A ‘\*’ in the plot for each cluster represents its cluster’s centroid. A well-calibrated teacher can effectively capture the inter-class relationships and serve as a reliable dynamic label smoothing prior such as MDCA [5]. Observe that the classes: **bowl** and **plate** are visually similar and hence the penultimate visualizations of these classes should be closer than the dissimilar class: **willow\_tree**. As the temperature  $T$  is increased the similar classes diffuse into one in the case of KD with LS while KD with MDCA offers better separation, retaining the semantic similarity while being well separated from the dissimilar class.

due to systematic diffusion among highly similar classes as seen in the penultimate representations. These discriminating features are crucial for achieving calibration by resolving confusion among similar classes. However, as  $T$  increases further, we simultaneously amplify the relationships between somewhat related classes [17], while diminishing the relationships between very similar classes. This nuanced understanding highlights the intricate interplay between temperature, class relationships, and calibration, shedding light on the optimal conditions for achieving calibration in KD scenarios.

#### 4 Illustration of the generality of KD(C) framework

Fig. S5 presents a novel framework KD(C) that leverages calibrated teachers through KD to produce DNNs with the least calibration error. This comprehensive framework encompasses the full spectrum, enabling models with varying capacity (smaller/larger) to distill student models with the least calibration error and better accuracy compared to the SOTA post-hoc/train-time calibration methods.





**Fig. S3: Visualization of penultimate layer’s activations** (Teacher = ResNet56, Student = ResNet8, Dataset = CIFAR100). We train ResNet8 using calibration techniques: KD with LS (Left column) and KD with MDCA (Right Column). We follow the same setup and procedure used in papers [ [12, 14]] We use two semantically similar classes (**man**, **boy**) and one semantically dissimilar class (**crab**). A ‘\*’ for each cluster represents its cluster’s centroid. A well-calibrated teacher can effectively capture the inter-class relationships and serve as a reliable dynamic label smoothing prior such as MDCA [5]. Observe that the classes: **man** and **boy** are visually similar and hence the penultimate visualizations of these classes should be closer than the dissimilar class: **crab**. As the temperature  $T$  is increased the similar classes diffuse into one in the case of KD with LS while KD with MDCA offers better separation, retaining the semantic similarity while being well separated from the dissimilar class.

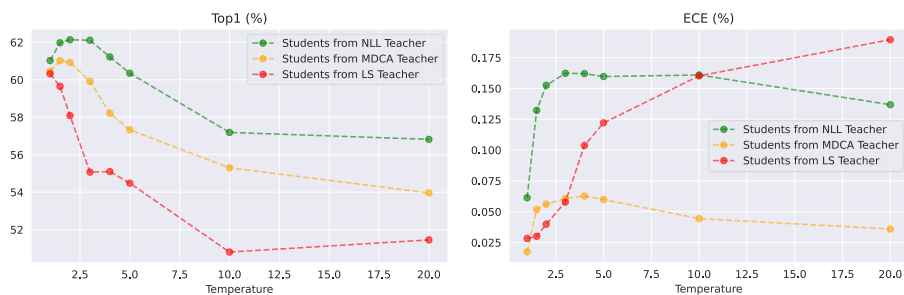
## 5 Additional Results

### 5.1 Reliability diagrams and Confidence Histograms

Reliability diagrams serve as effective visual aids for assessing the calibration of DNNs. They involve partitioning the predicted probabilities generated by DNNs into a predetermined number of bins along the  $x$ -axis. The  $y$ -axis represents the normalized count of events (e.g., class = “dog”) within each bin. A well-calibrated model will exhibit points that closely align with the main diagonal, spanning from the bottom left to the top right of the plot. Reliability diagrams corresponding to Fig. S6 are included to show that KD(C) variants obtain near SOTA results.

### 5.2 Effect of hyper-parameters like $T$ (temperature) and $\alpha$

We investigate the influence of hyperparameters  $T$  and  $\alpha$  on both calibrated and uncalibrated teacher models, as visually depicted in Fig. S7 (big-to-small) and Fig. S8 (small-to-big).



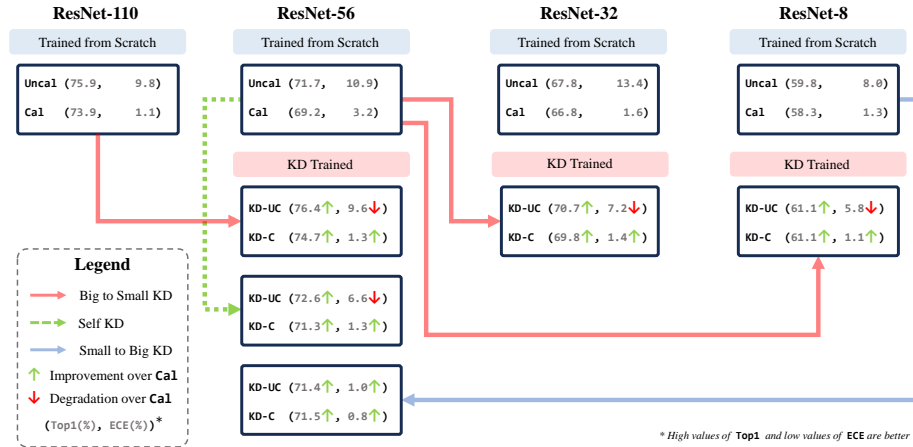
**Fig. S4: [Study of ECE variability in case of KD(C), specifically we consider KD with LS and KD with MDCA and study variation of accuracy and calibration as a function of temperature]:** Comparison of Top 1% accuracy and ECE when train-time calibration method is changed from Label Smoothing [16] and MDCA [5]: We use ResNet56 teacher on CIFAR100 and distill to ResNet8. Note that MDCA-based students have lower accuracy than NLL, however, ECE is largely stable when temperature  $T$  is varied.

**(a) Big teacher, Small student:** In this scenario as we increase the value of  $\alpha$ , we witness an intuitive rise in calibration. However, this effect is predominantly noticeable for small values of  $T$  (depicted in the bottom-right region of Fig. S7). Generally, the calibration errors (ECE) incurred by distilling students from a calibrated teacher tend to be markedly lower than those distilled from an uncalibrated teacher, as evident from the bottom row in Fig. S7. **(b) Small teacher, Big student:** Initially, we observe an expected trend: as  $\alpha$  increases (signifying a higher dependence on the teacher), accuracy experiences a decrease. This outcome arises from the process of distillation from a weaker teacher. However, when distilling from a calibrated teacher, we discern that elevating  $\alpha$  results in enhanced calibration. Nevertheless, this improvement in calibration is accompanied by a trade-off with accuracy.

Notably, we find that optimal calibration is generally achieved when  $T \approx 1$ , regardless of the size of the teacher model employed. This observation aligns with the findings presented in [15], which suggest that maximizing fidelity with the teacher model yields the best transfer of properties.

### 5.3 Calibration Performance under dataset drift

DNNs are found to be over-confident and highly uncalibrated under dataset/domain shift [19]. We investigate the robustness of our method KD(C) by examining the degradation in calibration under natural/non-semantic shift (images with the same label but different distribution). We carry out this study for ResNet56 pre-trained on the CIFAR100 dataset along with various calibration techniques and report the evaluation results on CIFAR100-C [6] in Fig. S9. We used ResNet56



**Fig. S5: An illustration of KD(C) framework’s generality using calibration method as MDCA.** We can distill a calibrated student from a large teacher and vice-versa yielding SOTA calibration without any trade-offs in accuracy. “Uncal” and “Cal” mean uncalibrated and calibrated teachers trained using NLL and a recent SOTA calibration technique [5] respectively. KD(UC) and KD(C) refer to students distilled using “Uncal” and “Cal” teachers respectively. Going from a large calibrated teacher to a smaller student yields SOTA calibrated student, with an additional boost in accuracy (E.g., compare ResNet56 “Trained from scratch” with ResNet56 “KD-trained” student from ResNet110). Self-distillation and going from a smaller teacher to a bigger student also have a similar effect on calibration, however, the gains in accuracy are comparable to respective models trained from scratch. The above results are on CIFAR100.

models that were trained with ResNet110 as teacher for KD(UC) and KD(C). We observe KD(UC) and KD(C) achieve the highest accuracy across all severities, with the latter achieving close to the best ECE (LS achieves best ECE), however, KD(C) achieves the best AUROC score in comparison to any other calibration technique. This indicates, KD(C) is better across all metrics measuring reliability (be it calibration or refinement, while also giving an additional boost in accuracy).

## 6 Results on ImageNet

Due to limited computational resources, conducting a wide range of experiments is challenging. To address this, we have employed ImageNet-100 (IN100 PyTorch Implementation: Training ResNets on ImageNet-100 at [tinyurl.com/mr3cr8rd](https://tinyurl.com/mr3cr8rd)) as a proxy to estimate performance on the full ImageNet-1K dataset. ImageNet-100 is a subset of ImageNet-1K, consisting of 100 randomly selected classes, allowing us to work with a higher resolution dataset ( $224 \times 224$ ) compared to the smaller Tiny ImageNet ( $64 \times 64$ ) dataset previously reported. ?? presents

**Table T1: [Self-distillation]** using MobileNetV2 feature extractor on Tiny-ImageNet dataset. Note that the main paper reported self-distillation results using the MobileNetV2 feature extractor on the CIFAR10 dataset. **Top3** best KD(C) variants are reported. KD with MDCA variant of **KD(C)** achieve competitive calibration results with SOTA.

Calibration Method	Top1 (%) ↑	ECE (%) ↓	SCE (%) ↓	ACE (%) ↓
NLL	50.43	13.72	0.24	13.72
LS [16]	51.20	3.84	0.19	3.88
CE with TS [4]	50.43	13.72	0.24	13.72
MMCE [9]	50.30	11.32	0.21	11.32
MixUp [18]	52.02	4.74	0.19	4.73
PSKD [7]	53.66	13.27	0.21	13.27
MDCA [5]	46.81	1.52	0.19	<b>1.11</b>
CPC [2]	51.27	12.01	0.21	12.01
MbLS [10]	50.11	8.87	0.20	8.87
ACLS [13]	<b>56.60</b>	7.13	<b>0.17</b>	7.06
KD with NLL	49.31	4.20	0.20	4.20
<b>Ours (KD with MDCA)</b>	45.79	<b>0.85</b>	0.21	<b>1.17</b>
<b>Ours (KD with LS)</b>	49.69	2.69	0.19	<b>2.68</b>
<b>Ours (KD with MbLS)</b>	49.33	2.86	0.20	2.89

Method	ImageNet-100					
	Top1 (%) ↑	AUROC ↑	ECE (%) ↓ $\times 10^{-2}$	SCE (%) ↓ $\times 10^{-3}$	ACE (%) ↓ $\times 10^{-3}$	smECE (%) ↓ $\times 10^{-2}$
NLL	67.56	99.68	05.37	00.31	<b>00.12</b>	05.37
<b>KD + NLL</b>	<b>67.67</b>	<b>99.69</b>	<b>04.79</b>	<b>00.29</b>	00.13	<b>04.67</b>
MDCA	<b>68.54</b>	99.62	07.01	00.32	00.16	07.01
<b>KD + MDCA (Ours)</b>	68.32	<b>99.64</b>	<b>02.84</b>	<b>00.27</b>	<b>00.13</b>	<b>02.82</b>
FL	67.90	98.84	67.70	<b>00.55</b>	01.74	38.81
<b>KD+FL+MDCA (Ours)</b>	<b>68.52</b>	<b>99.52</b>	<b>43.05</b>	00.65	<b>00.85</b>	<b>35.93</b>

**Table T1:** Illustration of generalisability to largescale datasets: ImageNet-1K

results for several methods, with all KD variants trained for only 25 epochs, yet still demonstrating improved calibration performance.

## 7 Computational Complexity

Post-hoc calibration methods are simple to implement and require only a small validation set with few parameters. In contrast, train-time methods involve all model parameters and do not need separate validation data, which can lead to better generalization and effectiveness. Training student models through KD from calibrated teachers adds minimal complexity compared to standard NLL training of students. However, KD(C) does require access to pre-calibrated teacher models. If such models are unavailable, KD(C) necessitates an initial step of training and calibrating a teacher, which introduces some computational overhead. Despite this potential additional step, the trade-off is generally worthwhile. The significant improvements in model reliability achieved through KD(C) justify the extra

computational effort, especially in applications where model trustworthiness is crucial.

**Table T1: [Effect of TS on KD(C)]:**The student is calibrated by distilling from an MDCA calibrated teacher KD with MDCA (a variant of KD(C)). The table shows that further temperature scaling (TS) does not impact the models trained with KD with MDCA as they are calibrated to start with. Parameters: WRN-40-1: 0.56M ; WRN-40-2: 2.24M ; MNV2: 2.25M

Dataset	Teacher Model	Student Model	Temperature	Top1 (%)	ECE (%)	SCE (%)	ACE (%)
CIFAR100	WRN-40-2	WRN-40-1	0.10	71.06	26.40	0.55	26.39
			0.20	71.06	23.79	0.52	23.79
			0.50	71.06	15.74	0.38	15.74
			0.75	71.07	8.44	0.27	8.44
			<b>1.00</b>	<b>71.06</b>	<b>0.98</b>	<b>0.20</b>	<b>1.10</b>
			1.25	71.06	8.60	0.27	8.60
			1.50	71.06	18.00	0.42	18.00
			1.75	71.06	27.11	0.58	27.11
			2.00	71.06	35.22	0.74	35.22
			2.25	71.06	41.99	0.88	41.99
			2.50	71.06	47.39	0.99	47.39
			2.75	71.06	51.60	1.07	51.60
			3.00	71.06	54.86	1.12	54.86
			3.25	71.06	57.37	1.13	57.37
			3.50	71.06	59.32	1.11	59.32
			3.75	71.06	60.86	1.07	60.86
			4.00	71.06	62.08	1.00	62.08
			4.25	71.06	63.06	0.92	63.06
	4.50	71.06	63.86	0.81	63.86		
	4.75	71.06	64.52	0.69	64.52		
	5.00	71.06	65.07	0.56	65.07		
	WRN-40-2	MNV2	0.10	68.67	30.64	0.64	30.64
			0.20	68.67	27.69	0.60	27.68
			0.50	68.67	18.50	0.44	18.50
			0.75	68.67	10.65	0.30	10.65
			<b>1.00</b>	<b>68.67</b>	<b>1.52</b>	<b>0.20</b>	<b>1.64</b>
			1.25	66.40	5.27	0.23	5.26
			1.50	68.67	13.09	0.34	13.09
			1.75	68.67	20.54	0.47	20.52
			2.00	68.67	27.34	0.60	27.34
			2.25	68.67	33.37	0.71	33.37
			2.50	68.67	38.53	0.81	38.53
			2.75	68.67	42.85	0.90	42.85
			3.00	68.67	46.39	0.96	46.39
3.25			68.67	49.27	1.00	49.27	
3.50			68.67	51.61	1.02	51.61	
3.75			68.67	53.50	1.01	53.50	
4.00	66.40	55.04	0.99	55.04			
4.25	68.67	56.31	0.95	56.31			
4.50	68.67	57.35	0.90	57.35			
4.75	68.67	58.22	0.83	58.22			
5.00	68.67	58.94	0.77	58.94			

## 8 Training details

In this section, we provide a detailed summary of the hyperparameters and training techniques used, in order to ensure reproducibility. All models have been trained on 40GB Nvidia A100 GPUs. The code was written using the PyTorch framework. We make use of automatic mixed precision training in order to reduce training time. We borrow some code from the official implementation of [5,11,20].

For CIFAR10/100 datasets, we train all `ResNets` / `WideResNets` models using a learning rate of 0.1 for 160 epochs. The learning rate is decayed by a factor of 10 at epoch 80 and 120. We use SGD optimizer with momentum 0.9 and weight decay of  $5e - 4$ . We use a batch size of 128. For the larger models like `ResNet-110`, we train them using a learning rate of 0.05 for 240 epochs. The learning rate is decayed by a factor of 10 at epoch 150, 180 and 210. We use SGD optimizer with momentum 0.9 and weight decay of  $5e - 4$ . We use a batch size of 64.

For `Tiny-ImageNet` dataset, all models are trained using a maximum learning rate of 0.1 with a cosine annealing learning rate with a warmup of 1000 steps with minimum learning rate  $1e - 5$ . The weight decay and momentum are  $5e - 4$  and 0.9 respectively. We train the models for 100 epochs with a batch size of 128.

For training students using KD, we use the same hyper-parameters for the respective datasets. For big-to-small KD (e.g. `WideResNet-40-2`  $\rightarrow$  `WideResNet-40-1`), we grid search  $T$  (temperature) and  $\alpha$  (distillation weight) in the ranges  $\{1, 1.5, 2, 3, 4, 5, 10, 20\}$  and  $\{0.9, 1.0\}$  respectively. For small-to-big KD and self-distillation (e.g. `MobileNetV2`  $\downarrow$  `DenseNet-121`, `MobileNetV2`  $\rightarrow$  `MobileNetV2`), we grid search  $T$  and  $\alpha$  in the ranges  $\{1, 1.5, 2, 3, 4\}$  and  $\{0.1, 0.3, 0.5, 0.7, 0.9, 1.0\}$  respectively.

For baselines, we use the recommended hyperparameters as suggested by the respective authors [2, 3, 5, 7, 8, 10, 13, 16, 18], i.e. for LS [16], we use smoothing of 0.1; for PSKD [7] we use  $\alpha = 0.8$ ; for MixUp [21] the mixup hyperparameter was taken as 0.4 as it was reported to be the best by the authors, for MDCA [5], we grid search for the best performing  $\beta \in \{1, 5, 10\}$  and  $\gamma \in \{1, 2, 3, 4, 5\}$ ; For MMCE [8], we grid search for the best performing  $\beta \in \{1, 2, 3, 4, 5\}$ . For ACLS [13], we set the margin (M) to 6 for CIFAR-10 and CIFAR-100, and 10 for `Tiny-ImageNet` dataset as recommended.

## 9 Reproducibility

In the spirit of reproducible research, we intend to make the source code available post-acceptance. To aid reviewers, the source code for our approach is attached along with the supplemental material. Details of our setup and implementation of the baselines can be found at: `Code/README.md` folder.

## 10 Limitations

While our work paves way to create optimal lightweight models that are both accurate and calibrated, it is important to acknowledge three potential limitations that we plan to address in future research - (a) principled approach to select hyperparameters, such as the temperature  $T$ , distillation weight  $\alpha$ , calibration regularization coefficient  $\beta$ , and characterization of optimal student-teacher capacity difference for best calibration, (b) extending theoretical insights to general nonlinear networks, (c) benchmarking KD(C) on natural language processing (NLP) tasks, particularly when the teacher networks belong to the family of large language models (LLMs).

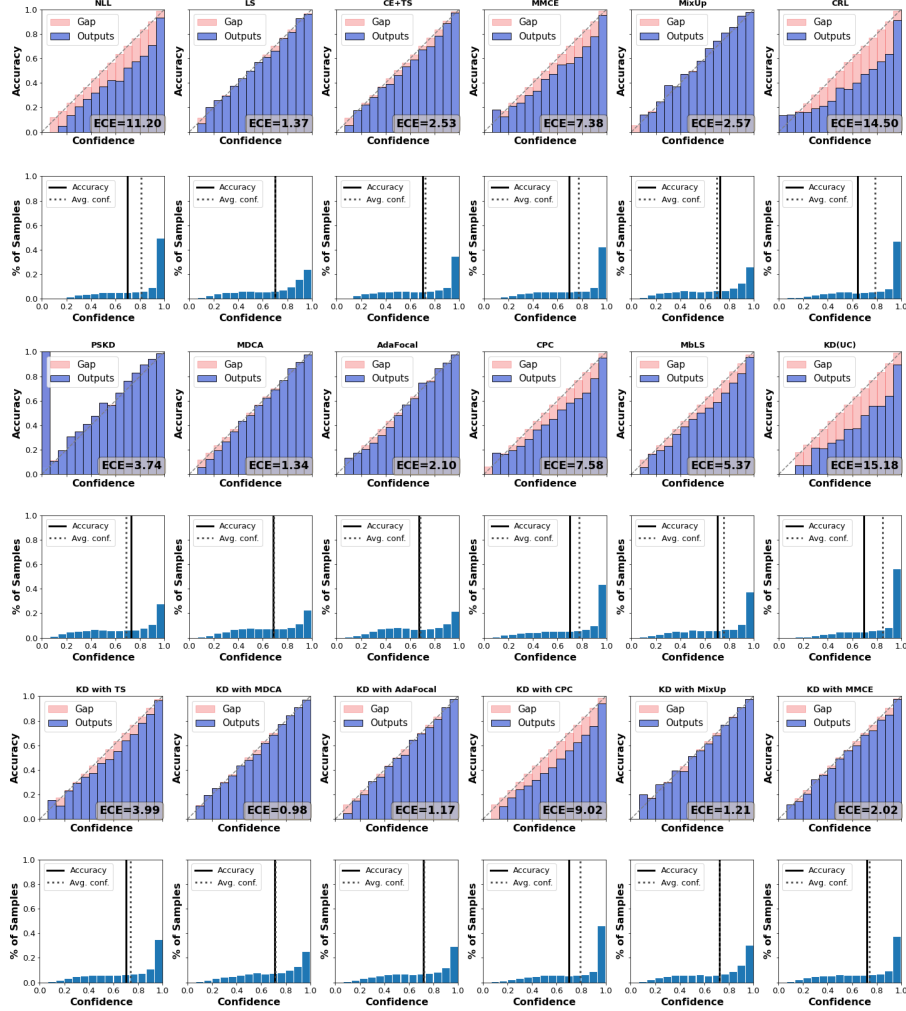
## 11 Broader Impact

Bigger DNN models aren't necessarily better models. From a deployment standpoint, the size of the weights affects the inference time and storage constraints on edge devices which is crucial in applications such as augmented reality and robotics. Our proposed algorithm has the potential to be employed in trustworthy lightweight models on the edge. In our endeavor to deploy lightweight models that are also reliable, we delve into the realm of knowledge distillation, extending its traditional function of transferring accuracy from teacher networks to student networks. Through this exploration, we have discovered a novel approach to calibrating models effectively. We present, arguably for the first time, compelling evidence that model calibration can be achieved without sacrificing accuracy through knowledge distillation. Notably, our implementation of knowledge distillation not only guarantees enhanced model calibration but also outperforms the accuracy obtained through conventional training from scratch in specific cases. This innovative approach enables us to simultaneously accomplish the dual objectives of optimal calibration and improved accuracy.

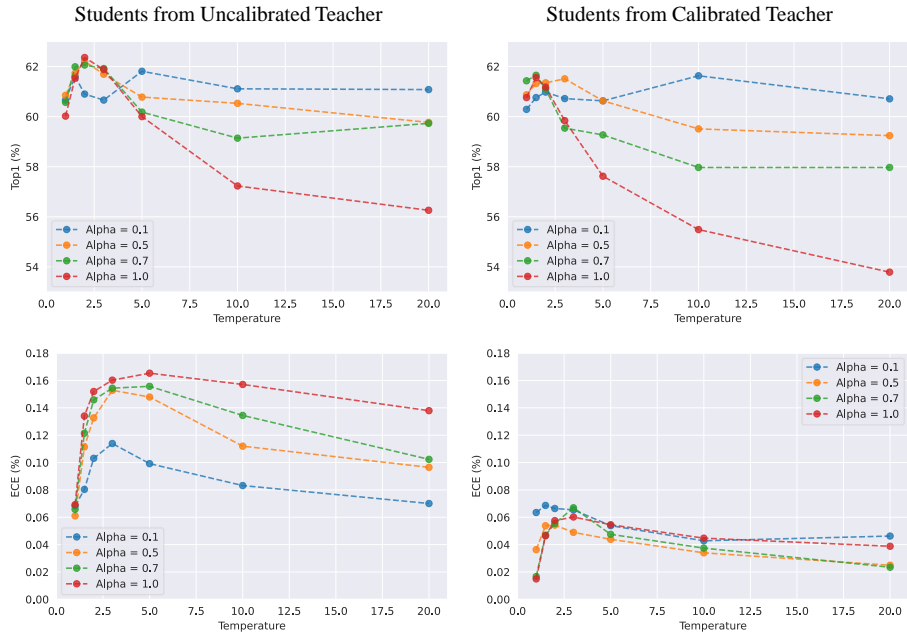
Towards this end, we provide extensive theoretical findings that extend beyond the realms of accuracy transfer and calibration alone. We show, through optics of linear teacher and student networks, that the optimization of student network weights through knowledge distillation enables them to exhibit similar behavior and performance as their respective teachers (see Theorem 1 in the main text). Subsequently, the scope of producing trustworthy models can also be extended to incorporate characteristics, such as fairness and refinement. On a more specific note, Theorem 2 in our work shows that there is a definite advantage of working with calibrated teachers over uncalibrated teachers, i.e., calibrated teachers tend to produce calibrated students without compromising on accuracy. Hence our approach, KD(C), centers around train-time calibration of teacher models, enabling them to generate accurate and optimally calibrated students through knowledge distillation. Significantly, based on our empirical evaluations, it is evident that the transfer of calibration operates bidirectionally. This means that larger calibrated models can be utilized to create smaller calibrated models, and conversely, smaller calibrated models can also serve as a foundation for generating larger calibrated models.

Overall, the research contributes to the advancement of model calibration, accuracy, trustworthiness, and scalability, which can have significant implications in various fields relying on the deployment of reliable and lightweight models.

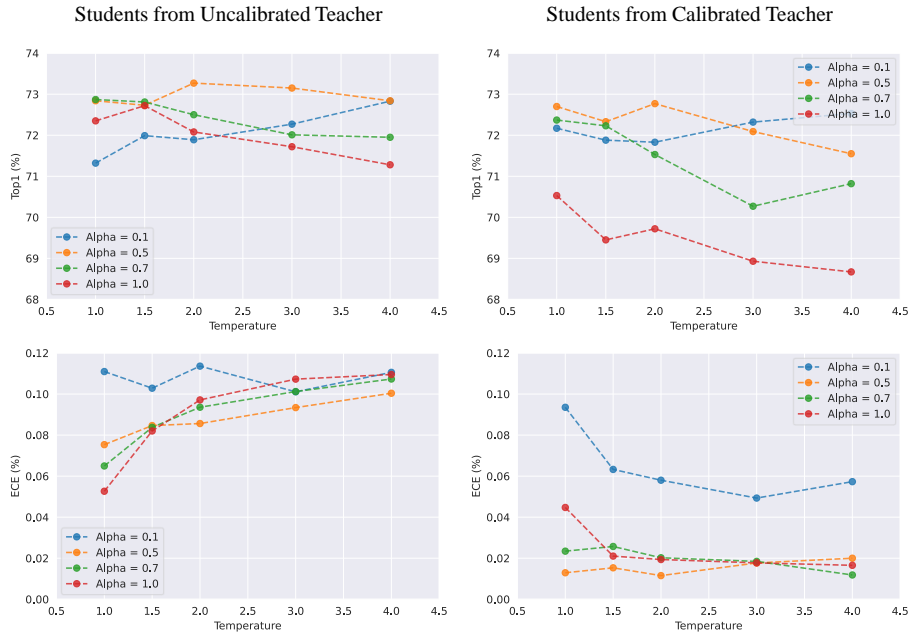




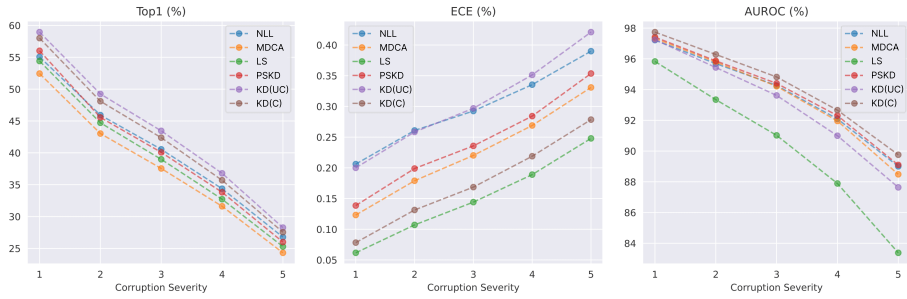
**Fig. S6:** Reliability Plots for top-5 KD with (Ours) techniques on WideResNet-40-1 on CIFAR100. Teacher used: WideResNet-40-2. KD(C) framework achieves competitive calibration results for KD with MDCA, KD with AdaFocal and KD with MixUp.



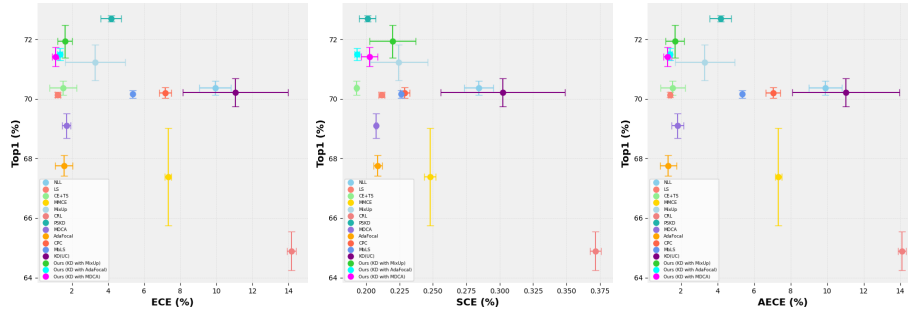
**Fig. S7:** We study the effect of varying temperature,  $T$  and distillation weight  $\alpha$ , on ECE and top 1% accuracy when ResNet56 teacher model is used and ResNet8 as student on CIFAR100 dataset. Observe the optimal values of ECE and top 1% accuracy when  $T$  is set around 1. For calibration KD with MDCA was used.



**Fig. S8:** We study the effect of varying temperature,  $T$  and distillation weight  $\alpha$ , on ECE and top 1% accuracy when ResNet32 teacher model is used and ResNet56 as student on CIFAR100 dataset. KD with MDCA was used for calibration.



**Fig. S9:** Robustness to corruption, tested on CIFAR100-C dataset [6] using ResNet-56. KD(UC) and KD(C) were trained using ResNet-110 as a Teacher. Note that KD(C) provides a good trade-off between accuracy and calibration, at the same time achieving the highest AUROC (even though LS outperforms KD(C) by a tiny margin in terms of calibration, KD(C) has significantly better AUROC and accuracy. AUROC indicates better inter-class separability in classifiers thereby enhancing trustworthiness in addition to calibration benefits. KD(C) uses KD with MDCA variant.



**Fig.S7: Comparative study of accuracy vs. calibration trade-offs associated with existing calibration techniques and ours (Top-left is most preferred):** The mean and one standard scatter error bars for Top1, ECE and SCE of WideResNet-40-1 trained on CIFAR100 using SOTA calibration techniques. WideResNet-40-2 was used as Teacher for KD(UC) and the proposed, KD(C) variants. Note: KD(C) variants (magenta, cyan, and green) achieve the best results in terms of ECE, ACE and SCE, along with slight boosts in Top1 (an inherent KD-property). Further, the lower variances emphasize the reliability of KD(C) variants. All plots were generated by training WideResNet-40-1 models through every calibration technique on 3 runs.

## References

1. Chandrasegaran, K., Tran, N.T., Zhao, Y., Cheung, N.M.: Revisiting label smoothing and knowledge distillation compatibility: What was missing? In: International Conference on Machine Learning. pp. 2890–2916. PMLR (2022) [7](#)
2. Cheng, J., Vasconcelos, N.: Calibrating deep neural networks by pairwise constraints. In: 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 13699–13708 (2022). <https://doi.org/10.1109/CVPR52688.2022.01334> [7](#), [12](#), [14](#)
3. Ghosh, A., Schaaf, T., Gormley, M.: Adafocal: Calibration-aware adaptive focal loss. In: Advances in Neural Information Processing Systems. vol. 35, pp. 1583–1595 (2022), [https://proceedings.neurips.cc/paper\\_files/paper/2022/file/0a692a24dbc744fca340b9ba33bc6522-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/0a692a24dbc744fca340b9ba33bc6522-Paper-Conference.pdf) [14](#)
4. Guo, C., Pleiss, G., Sun, Y., Weinberger, K.Q.: On calibration of modern neural networks. CoRR **abs/1706.04599** (2017), <http://arxiv.org/abs/1706.04599> [12](#)
5. Hebbalaguppe, R., Prakash, J., Madan, N., Arora, C.: A stitch in time saves nine: A train-time regularizing loss for improved neural network calibration. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 16081–16090 (June 2022) [1](#), [7](#), [8](#), [9](#), [10](#), [11](#), [12](#), [14](#)
6. Hendrycks, D., Gimpel, K.: A baseline for detecting misclassified and out-of-distribution examples in neural networks. CoRR **abs/1610.02136** (2016), <http://arxiv.org/abs/1610.02136> [10](#), [19](#)
7. Kim, K., Ji, B., Yoon, D., Hwang, S.: Self-knowledge distillation with progressive refinement of targets. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). pp. 6567–6576 (October 2021) [12](#), [14](#)
8. Kumar, A., Liang, P., Ma, T.: Verified uncertainty calibration. arXiv preprint arXiv:1909.10155 (2019) [14](#)
9. Kumar, A., Sarawagi, S., Jain, U.: Trainable calibration measures for neural networks from kernel mean embeddings. In: Dy, J., Krause, A. (eds.) Proceedings of the 35th International Conference on Machine Learning. Proceedings of Machine Learning Research, vol. 80, pp. 2805–2814. PMLR (10–15 Jul 2018), <http://proceedings.mlr.press/v80/kumar18a.html> [12](#)
10. Liu, B., Ben Ayed, I., Galdran, A., Dolz, J.: The devil is in the margin: Margin-based label smoothing for network calibration. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 80–88 (2022) [7](#), [12](#), [14](#)
11. Mukhoti, J., Kulharia, V., Sanyal, A., Golodetz, F., Torr, P., Dokania, P.: Calibrating deep neural networks using focal loss. Advances in Neural Information Processing Systems **33**, 15288–15299 (2020) [14](#)
12. Müller, R., Kornblith, S., Hinton, G.E.: When does label smoothing help? Advances in neural information processing systems **32** (2019) [7](#), [8](#), [9](#)
13. Park, H., Noh, J., Oh, Y., Baek, D., Ham, B.: Acls: Adaptive and conditional label smoothing for network calibration. In: Proceedings of the IEEE/CVF ICCV (2023) [7](#), [12](#), [14](#)
14. Shen, Z., Liu, Z., Xu, D., Chen, Z., Cheng, K.T., Savvides, M.: Is label smoothing truly incompatible with knowledge distillation: An empirical study. arXiv preprint arXiv:2104.00676 (2021) [7](#), [8](#), [9](#)
15. Stanton, S., Izmailov, P., Kirichenko, P., Alemi, A.A., Wilson, A.G.: Does knowledge distillation really work? Advances in Neural Information Processing Systems **34**, 6906–6919 (2021) [10](#)

16. Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., Wojna, Z.: Rethinking the inception architecture for computer vision. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2015) 7, 8, 10, 12, 14
17. Tang, J., Shivanna, R., Zhao, Z., Lin, D., Singh, A., Chi, E.H., Jain, S.: Understanding and improving knowledge distillation. arXiv preprint arXiv:2002.03532 (2020) 8
18. Thulasidasan, S., Chennupati, G., Bilmes, J., Bhattacharya, T., Michalak, S.: On mixup training: Improved calibration and predictive uncertainty for deep neural networks. In: Advances in neural information processing systems (2019) 12, 14
19. Tomani, C., Gruber, S., Erdem, M.E., Cremers, D., Buettner, F.: Post-hoc uncertainty calibration for domain drift scenarios. CoRR **abs/2012.10988** (2020), <https://arxiv.org/abs/2012.10988> 10
20. Yuan, L., Tay, F.E.H., Li, G., Wang, T., Feng, J.: Revisiting knowledge distillation via label smoothing regularization. In: CVPR (2021) 14
21. Zhang, H., Cissé, M., Dauphin, Y.N., Lopez-Paz, D.: mixup: Beyond empirical risk minimization. In: ICLR (2018) 14