

# SOCRATES LOSS FOR TRAINING AD-HOC CALIBRATED SELECTIVE CLASSIFIERS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Model reliability is paramount for critical real-world applications. To enhance reliability, it is essential to quantify uncertainty in model predictions, as achieved through Confidence Calibration and Selective Classification. Confidence Calibration ensures prediction confidences accurately reflect the actual likelihood of correctness, while Selective Classification allows a model to abstain from making predictions when uncertain. Although related, existing methods address each aspect separately, or both through post-hoc methods. Only one method, Confidence-aware Contrastive Learning for Selective Classification (CCL-SC), combines both in an ad-hoc manner. Despite being a powerful calibrator, CCL-SC has some drawbacks, including the absence of an additional unknown class, the use of two different losses (detrimental for calibration), and its cumbersome implementation. In the pursuit of reliable models and motivated by the idea of creating an ad-hoc calibrated selective classifier with an unknown class, we first empirically analyze the Self-Adaptive Training (SAT) method, a leading method in ad-hoc selective classification. We identify that while SAT excels in selective classification, it falls short in confidence calibration, especially when training for a small number of epochs (*e.g.*,  $\leq 100$ ). To address this, we introduce an original method that uses an unknown class and a unique novel loss, *Socrates loss*, which serves as a classifier and a calibrator with a unified optimization goal. This method mitigates overfitting and ensures theoretically well-calibrated predictions across all epochs, addressing the drawbacks of both CCL-SC and SAT, without the need for post-hoc processing or additional data. We integrate our method into the SAT implementation and extend it to provide selective classification and confidence calibration metrics. We show empirically that our method matches or improves the selective classification error rate of SAT and CCL-SC, while producing well-calibrated models in an ad-hoc manner through the evaluation on 6 image benchmark datasets across two architectures, VGG-16 and ResNet-34.

## 1 INTRODUCTION

Reliability, the ability of a model to consistently operate in real-world environments (Tran et al., 2022), becomes particularly important in critical real-world scenarios, including but not limited to medical diagnosis (Gireesh & Gurupur, 2023), nuclear security (Ayodeji et al., 2022), and biosecurity (McEwen et al., 2021). A reliable model should not only achieve strong predictive performance but also excel in the representation of its own uncertainty. To quantify uncertainty, different methods measure distinct aspects of the predictive uncertainty stemming from reliability, such as Confidence Calibration and Selective Classification. Selective classification allows models to abstain from making predictions when uncertain, ensuring cautious decision-making in high-risk applications. On the other hand, confidence calibration ensures that predictive confidence accurately reflects the likelihood of correctness. Although both strategies aim to enhance reliability, they are typically approached and studied independently (Zhang et al., 2023).

In critical high-risk applications, where trust in predictions is essential, integrating confidence calibration with selective classification is crucial. Recent work has highlighted this need and proposed new post-hoc methods (Fisch et al., 2022; Galil et al., 2023; Moon et al., 2020), and, to the extent of our knowledge, only one ad-hoc method, Confidence-aware Contrastive Learning for Selective Classification (CCL-SC) (Wu et al., 2024). Despite the fact that CCL-SC is able to output cal-

054 ibrated models, it has several drawbacks. Firstly, following the work of Feng et al. (2023), the  
055 extra unknown class was not added. The use of an extra unknown class is related to adaptations  
056 aiming to address Open-Set Recognition (OSR) in deep neural networks (Bendale & Boult, 2016;  
057 Patrick Schlachter & Yang, 2019). In the search for reliable models, we argue that mechanisms for  
058 OSR should be integrated with calibration and selective classification to enhance model reliability  
059 and adaptability. A reliable model should not only adapt its predictions to new scenarios but also be  
060 flexible enough to handle different use cases. Incorporating an unknown class provides more flexi-  
061 bility, allowing the model to function as a selective classifier or a traditional classifier with or without  
062 an extra unknown class. Secondly, the CCL-SC method features two losses: Confidence-aware Su-  
063 pervised Contrastive (CSC) loss for calibration and cross-entropy (c.e.) loss for classification. We  
064 have identified that this method miscalibrates the model once it reaches a certain calibration point  
065 which is due to the c.e. effect (see Section 4.1.2). Therefore, having two losses can be detrimental  
066 if one of the losses is not specifically focused on calibration, highlighting the need for a unified loss  
067 with the same optimization goal. Thirdly, CCL-SC exhibits variable behavior in terms of Expected  
068 Calibration Error (ECE), depending on the architecture and dataset. The loss function across epochs  
069 also shows varied trends, including spikes, which could suggest training instability. Finally, the  
070 implementation of the CCL-SC method requires extensive modifications to the training code.

071 Motivated by these drawback and the idea of creating an ad-hoc calibrated selective classifier with  
072 a capability to estimate the probability for an unknown class, we first empirically analyzed the  
073 calibration capability of the Self-Adaptive Training (SAT) (Huang et al., 2020), the state-of-the-art  
074 for end-to-end selective classification with an extra unknown class. This analysis showed that SAT  
075 loss does not seem to be a calibration loss, as in the case of training for a smaller number of epochs  
076 (*e.g.*,  $\leq 100$ ) or using hard-to-classify datasets like CIFAR-100 and Food-101.

077 We propose a method to train calibrated selective classifiers by introducing an extra unknown class  
078 and using a novel unified loss, *Socrates* loss. *Socrates* loss owes its name to the famous quote  
079 of the philosopher Socrates: *I know that I know nothing*; which reflects the power of the loss to  
080 train a model to be aware of its own uncertainty. This loss integrates classification and calibration  
081 into a single optimization problem, optimizing a single loss function, and does need several losses,  
082 switching to a different loss during training, or post-hoc processing. The loss uses its knowledge  
083 about when it does not know, and dynamically utilizes model predictions to guide training, by giving  
084 more attention to hard-to-classify instances.

085 We empirically evaluated our method to SAT and CCL-SC on the CIFAR-10, CIFAR-10C, CIFAR-  
086 100, CIFAR-100C, SVHN, and Food-101 datasets with VGG-16 and ResNet-34 architectures. In  
087 terms of calibration across epochs, our method outperforms SAT and is comparable to CCL-SC,  
088 while effectively addressing the previously discussed drawbacks of both methods. Moreover, our  
089 method achieves similar or lower Selective Classification error rates compared to CCL-SC and SAT,  
090 with notable improvements observed over SAT on the hard-to-classify CIFAR-100, CIFAR-100C  
091 and Food101 datasets. For instance, using the Food-101 dataset, the VGG-16 architecture, and  
092 100% of coverage, we achieve Selective Classification error rates (%) of 26.93 for *Socrates*, 68.23  
093 for SAT and 27.18 for CCL-SC.

094 To summarise, our contributions are as follows:

- 095 • An easy-to-implement ad-hoc method that uses an extra unknown class and a novel loss,  
096 *Socrates* loss, integrating classification and calibration into a unified optimization goal.
- 097 • A Python implementation to train *Socrates*, SAT and CCS-CL, and evaluate selective clas-  
098 sification and confidence calibration. In addition, the code to reproduce the results of this  
099 paper is also provided.
- 100 • A theoretical analysis that proves *Socrates* loss a) forms a regularize upper bound in the  
101 Kullback-Leibler divergence, avoiding overconfident predictions and improving calibra-  
102 tion. b) acts as a regularizer (of the network weights) when the model is sufficiently confi-  
103 dent, avoiding miscalibration and overfitting.
- 104 • A comparative empirical analysis of *Socrates*, SAT, CCL-CS, in terms of calibration and  
105 selective classification performance, on 6 benchmark datasets across two network architec-  
106 tures.  
107

## 2 RELATED WORK

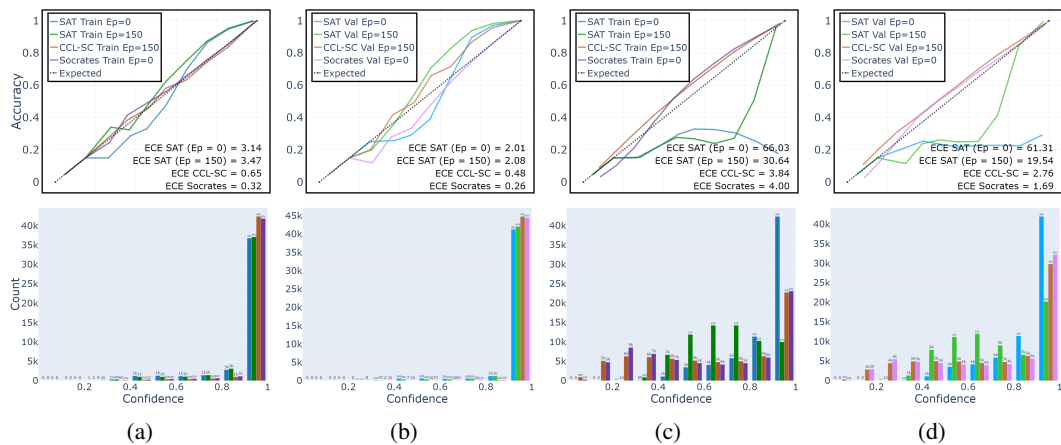


Figure 1: Reliability Diagrams of the last epoch displaying ECE values for CIFAR-10 (Figures 1a and 1b) and Food-101 (Figures 1c and 1d) datasets trained with the VGG-16 architecture using SAT, CCL-SC and Socrates methods.

In 2017, Guo et al. (2017) revealed that modern Neural Networks (NNs) are no longer well-calibrated and exhibit overconfidence. The degree of confidence calibration in NNs can be illustrated through visual representations and quantifiers (Section 3.2). To address the issue of miscalibration, the research community has focused on developing new post-hoc and ad-hoc methods. Post-hoc methods are applied to a trained model in a post-training process, such as Platt Scaling (Platt, 2000) and Temperature Scaling (Guo et al., 2017). Ad-hoc methods enhance both accuracy and calibration during training, creating end-to-end compact models by incorporating explicitly or implicitly a secondary optimization objective related to the predictive uncertainty of the model in the training objective (Liu et al., 2023). Although they remain underexplored, one way to achieve calibration is through the use of a specific loss function, such as Focal loss (Lin et al., 2020; Mukhoti et al., 2020). According to Zhang et al. (2023), despite the effectiveness of post-hoc methods, future models should integrate calibration into the training process. To that end, we focus our method on calibrating through a loss function.

Alternatively, one can train reliable models by considering the option to reject a prediction when the model is uncertain. Selective Classification (Geifman & El-Yaniv, 2017) can be addressed by post-hoc and ad-hoc methods. Post-hoc methods perform selective classification after training, such as LeCun et al. (1989) and Geifman & El-Yaniv (2017). Ad-hoc methods change the NNs training process and add extra heads or logits. Feng et al. (2023) divide these methods into learn to select, such as SelectiveNet (Geifman & El-Yaniv, 2019), and learn to abstain, such as Self-Adaptive Training (SAT) (Huang et al., 2020), methods. According to Feng et al. (2023), adding an extra head/logit is unnecessary for the ad-hoc Selective Classification problem, and the Softmax Response is the only selection mechanism required. We argue the opposite, that in the search for reliable models, the ability to handle unknown classes is valuable (Subsection 6.2).

Whereas some authors have emphasized the necessity of having calibrated selective classifiers (Zhang et al., 2023), the majority of published methods (Fisch et al., 2022; Galil et al., 2023; Moon et al., 2020) have focused on post-hoc integration, which presents several drawbacks: often requiring additional data, increasing the risk of bias, diffusing the optimization goal, and sometimes failing to fit when the calibration error is too complex. According to Zhang et al. (2023), a well-calibrated model could not be a good discriminator and vice versa. Although SAT was designed as a loss function to prevent overfitting, it has not been proven to be a promising ad-hoc calibration loss function across all epochs. Wu et al. (2024) presented CCL-SC, which is currently the state-of-the-art for ad-hoc calibrated selective classifiers. In contrast, Fisch et al. (2022) proposed a selective classifier that rejects instances based on calibration rather than potential misclassification, which represents a different goal from ours.

### 3 PROBLEM FORMULATION

We frame the problem as a multi-class classification task with  $(c + 1)$  classes, where the last class represents the unknown class.

#### 3.1 PROBLEM SETTING: SELECTIVE CLASSIFICATION

Selective Classification trades classifier coverage against accuracy. It is the ability of a model to reject instances when there is uncertainty. The rejected instances are potential out-of-distribution or lie in the tail of the data distribution; making predictions only on samples with confidence.

Let  $\mathcal{X}$  be the feature space,  $\mathcal{Y}$  be the label space, and  $P(\mathcal{X}, \mathcal{Y})$  be the data distribution over  $\mathcal{X} \times \mathcal{Y}$ . A selective model is a pair  $(f, g)$ , where the prediction function is  $f : \mathcal{X} \rightarrow \mathcal{Y}$ , in our case a classifier, and the selection function is  $g : \mathcal{X} \rightarrow \{0, 1\}$ . Then,  $f(x)$  makes a prediction when  $g(x) = 1$ , and abstains from making a prediction when  $g(x) = 0$ .

The performance of a selective classifier can be evaluated in terms of cost-sensitive learning (Cortes et al., 2016) where the rejection cost needs to be specified, or from a Risk-Coverage perspective (El-Yaniv & Wiener, 2010). Since specifying costs can be challenging (Geifman & El-Yaniv, 2017), we evaluate our method using the Risk-Coverage perspective. Coverage is defined as the probability mass of the nonrejected region of  $\mathcal{X}$ ,  $\phi(g) = \mathbb{E}[g(X)]$ . In practice, a soft selection function  $\tilde{g} : \mathcal{X} \rightarrow \mathbb{R}$  is often used, constraining the coverage with a threshold  $\tau \in \mathbb{R}$ . Then  $g$  is defined as  $g(x) := \mathbf{1}\{\tilde{g}(x) \geq \tau\}$ . Given a loss function, the selective risk, which corresponds to the selective error when the loss is 0/1, with respect to  $P$  can be defined as  $R(f, g) = \mathbb{E}[\mathcal{L}(f(X), Y) | g(X) = 1] = \frac{\mathbb{E}[\mathcal{L}(f(X), Y) | g(X)]}{\phi(g)}$ . This shows a dependency between risk and coverage; rejecting samples results in lower selective risk and lower coverage. Therefore, from a Risk-Coverage perspective, the minimization problem given a target coverage is:  $\min R(f, g)$  s.t.  $\phi \geq c_{\text{target}}$ .

We follow the Selective Classification problem for ad-hoc methods to train end-to-end selective classifiers proposed by SAT (Huang et al., 2020) and DeepGamblers (Liu et al., 2019), where the selection function  $g(\cdot)$  is replaced by  $f(\cdot)_c$  where  $c$  is the number of classes. In our proposed method, similar to SAT and DeepGamblers, an additional unknown class  $(c + 1)$  represents abstention.

#### 3.2 PROBLEM SETTING: CONFIDENCE CALIBRATION

Confidence calibration is the process of aligning predictive confidence with the actual likelihood of correctness, i.e. accuracy in the multiclass case. One method to reach ad-hoc calibration is through a loss function as with Focal loss (Lin et al., 2020). The confidence calibration level of a NN can be represented through visualizations and quantifiers.

A popular method for visualising confidence calibration is the Reliability Diagrams (Niculescu-Mizil & Caruana, 2005), which plot the expected sample accuracy as a function of confidence. Confidences can be grouped in different forms (Filho et al., 2023; Guo et al., 2017; Nguyen & O’Connor, 2015) to estimate expected accuracy from finite samples. In this work, we adopt the approach of Guo et al. (2017), grouping confidences into  $M$  interval bins of size  $1/M$ , increasing the probability of having multiple samples per estimation range. Let  $B_m$  be the test set of indices of samples whose confidence falls into the  $m$ -th bin,  $I_m = (\frac{m-1}{M}, \frac{m}{M}]$ . The confidence of bin  $B_m$  is estimated as  $\text{conf}(B_m) = \frac{1}{|B_m|} \sum_{i \in B_m} \hat{p}_i$ ; where  $\hat{p}_i$  is the confidence for sample  $i$ . The average accuracy is estimated as  $\text{acc}(B_m) = \frac{1}{|B_m|} \sum_{i \in B_m} \mathbf{1}(\hat{y}_i = y_i)$ ; where  $\hat{y}_i$  is the predicted class label and  $y_i$  is the true class label for sample  $i$ .

To measure calibration the most popular metrics are the Expected Calibration Error (ECE) and the Maximum Calibration Error (MCE) (Naeini et al., 2015). ECE is the weighted average of the difference between accuracy and confidence in each bin:  $ECE = \sum_{m=1}^M \frac{|B_m|}{n} |\text{acc}(B_m) - \text{conf}(B_m)|$ . MCE is the worst-case deviation and is valuable for high-risk frameworks:  $MCE = \max_{m \in \{1, \dots, M\}} |\text{acc}(B_m) - \text{conf}(B_m)|$ . It is common to use Brier Score as in Fisch et al. (2022) but, as it is an aggregate measure (Hernández-Orallo et al., 2012), it is inadequate for analyzing calibration in isolation.

An example of Reliability Diagrams along with the ECE values is presented in Figure 1.

## 4 OUR METHOD: CALIBRATED SELECTIVE CLASSIFICATION WITH AN UNKNOWN CLASS

We propose a versatile method that can be used as a selective classifier or as a standard classifier with or without an unknown class; well-calibrated in all cases. Inspired by the calibration principles of Focal loss and influenced by the selective classification power of SAT, we introduce a method to train calibrated selective classifiers by integrating an additional unknown class, referred to as *idk*, and using an easy-to-implement novel loss called Socrates loss, which maintains a unified optimization objective of classification and calibration.

Therefore, a classifier  $f(\cdot)_c$  is optimized by minimizing Socrates loss, which is defined as:

$$\mathcal{L}(f) = -\frac{1}{n} \sum_{i=1}^n (1 - \hat{p}_{i,y_i})^\gamma [t_{i,y_i} \log \hat{p}_{i,y_i} + \alpha_{dynamic} (1 - t_{i,y_i}) \log \hat{p}_{i,idk}]. \quad (1)$$

where  $\hat{p}_{i,y_i}$  is the prediction associated with the ground truth class and  $\hat{p}_{i,idk}$  is the prediction associated with the *idk* class,  $n$  the number of instances,  $\gamma$  a modularity factor controlling the down-weighting of easy examples (higher factor gives more weight to hard-examples), and  $\alpha_{dynamic}$  is a regularizer which controls attention to the unknown knowledge.

Initially, during the first selected  $E_s$  training epochs, the target is the ground truth label,  $t_i \leftarrow y_i$ . After  $E_s$ , the target is updated in each epoch as  $t_i \leftarrow \alpha_{momentum} \times t_i + (1 - \alpha_{momentum}) \times \hat{p}_{i,y_i}$  s.t.  $\alpha_{momentum} \in (0, 1]$ . This dynamic behaviour balances the importance of current predictions associated with the ground truth class and the *idk* class, reducing prediction instability. Our main proposal uses  $E_s = 0$ , thereby creating an end-to-end loss.

The logic behind the loss can be described as follows. If a sample was previously predicted with high confidence, the first part of the equation has more influence, resembling Focal loss and giving a bigger penalty towards hard-to-classify samples. This method helps to avoid overfitting and calibrate the model when the uncertainty is low. Conversely, if the sample seems uncertain (i.e., low previous confidence), the second part of the equation assumes greater importance acting as a selection function in the selective classifier. This part is influenced by an  $\alpha_{dynamic} \leftarrow (\max_{y_i \neq y_{gt}} \hat{p}_{i,y_i}) - \hat{p}_{i,y_{idk}}$ , which adjusts attention based on the awareness of the classifier of its own uncertainty, of its own lack of knowledge. If the classifier recognizes its own uncertainty, i.e., the *idk* class predicted probability surpasses other class probabilities (without the ground truth class probability), only the first part is considered; as the model knows it does not know. Otherwise, if the classifier is not aware of its lack of knowledge, the selection function gains relevance weighted by the focal component. This method increases penalties for hard-to-classify instances and for instances where the classifier does not have certainty that it does not know.

The pseudocode of the method can be found in Appendix A and a mathematical example in B.

### 4.1 THEORETICAL ANALYSIS

#### 4.1.1 SOCRATES LOSS FORMS A REGULARIZED UPPER BOUND IN THE KULLBACK-LEIBLER DIVERGENCE

It is well-known that c.e. loss minimizes (provides an upper bound for) the Kullback-Leibler (KL) divergence between the predicted and the target distributions over classes, i.e.,  $\mathcal{L}_{c.e.}(f) \geq D_{KL}(q||\hat{p})$ . KL divergence quantifies the information difference between two distributions. In our case, Socrates loss minimizes KL divergence while regularizing by increasing the entropy of the predicted distribution and leveraging the predictions associated with the unknown class. The regularization parameters are  $\gamma$ ,  $\alpha_{dynamic}$ , and  $\Delta_{reg}$ ; where  $\Delta_{reg} = (1 - t_y)[\gamma \hat{p}_y \log \hat{p}_{idk} - \log \hat{p}_{idk}]$ . Therefore:

$$\mathcal{L}(f) \geq D_{KL}(q||\hat{p}) - \gamma \mathbb{H}[\hat{p}] + \alpha_{dynamic} \Delta_{reg}. \quad (2)$$

This regularised entropy increase, along with the regularization applied through the prediction associated with the unknown class, prevents the model from becoming overconfident. Then, substituting

the c.e. loss with Socrates loss incorporates a maximum-entropy regulariser (Pereyra et al., 2017) to the KL minimization. As demonstrated by Lin et al. (2020), higher entropy can prevent overconfident predictions, improving model calibration. Therefore, Socrates Loss forms a regularize upper bound in the KL divergence, avoiding overconfident predictions and improving calibration. The proof can be found in Appendix C.1.

#### 4.1.2 SOCRATES LOSS REGULARIZES THE WEIGHTS OF THE NETWORK

Guo et al. (2017) and Lin et al. (2020) proved there is a relationship between miscalibration and overfitting (but not the opposite). This occurs when the loss function attempts to further reduce its value even after perfect high confidence has been achieved. Lin et al. (2020) demonstrated that for misclassified samples using c.e. loss the network progressively grows more confident in its incorrect predictions. Socrates loss acts as a regularizer with an increased penalty highly associated with the unknown class when the model begins to overfit. Furthermore, the norms of the weights are higher at the beginning of the training compared to those trained with c.e. It is when the model starts being miscalibrated that there is a change in the ordering of the weight norms, due to a big increase in the weight norm of the models with c.e. This behaviour shows that Socrates loss acts as a regularizer when the model is sufficiently confident, avoiding miscalibration and overfitting.

Therefore, let  $\mathcal{L}_{c.e.}(f)$  be c.e. loss, and  $\mathcal{L}(f)$  be Socrates loss. The gradients of the neural network trained with  $\mathcal{L}(f)$  are smaller than the ones trained with  $\mathcal{L}_{c.e.}(f)$  when a perfect confidence is reached and the model could start overfitting and then become miscalibrated, i.e.,

$$\left\| \frac{\partial \mathcal{L}(f)}{\partial w} \right\| \leq \left\| \frac{\partial \mathcal{L}_{c.e.}(f)}{\partial w} \right\|. \quad (3)$$

The proof can be found in Appendix C.2.

## 5 EXPERIMENT SETTINGS

For the upcoming experiments, we initially evaluated SAT against Focal loss, a calibration loss function. Afterwards, we evaluated the Socrates method, comparing it to the SAT and CCL-SC methods. To this end, we extended the publicly available SAT implementation to create a framework for training and evaluating calibrated selective classifiers.<sup>1</sup>

Table 1: Specifications of datasets employed in the experimental phase.

Dataset	Image Size	Classes	Train	Test	Specifications
CIFAR-10	32x32x3	10	50000	10000	Easy-to-classify
CIFAR-10C	32x32x3	10		50000	10000 using 5 levels of corruption
CIFAR-100	32x32x3	100	50000	10000	Hard-to-classify dataset
CIFAR-100C	32x32x3	100		50000	10000 using 5 levels of corruption
SVHN	32x32x3	10	73257	26032	Easy-to-classify real-world dataset
Food-101	224x224x3	101	75750	25250	Hard-to-classify real-world dataset

We used a VGG-16 and a ResNet-34 architecture for the datasets specified in Table 1. Each configuration was trained with five different seeds. Additional hyper-parameters details are in Appendix D.

SAT, CCL-SC and Socrates methods can be initialized in the first epochs with another loss (e.g., c.e. or Focal loss), and then switched to the main loss. For Selective Classification, the SAT authors instantiated the number of first epochs at 0, and for CCL-SC at 150. Therefore, we conducted two experiments: *first-epochs* ( $E_s = 150$ ) for Socrates and SAT methods, and *end-to-end* ( $E_s = 0$ ) for Socrates with Focal and Socrates losses, SAT with c.e. and SAT losses, and CCL-SC with c.e. and CSC losses. The goal is to determine whether a method with a unified loss, i.e., *first-epochs* case with Socrates method using only the Socrates loss, can achieve or surpass similar selective classification results while addressing the calibration issues of SAT and CCL-SC.

<sup>1</sup>The code is publicly available at <https://anonymous.4open.science/r/Socrates>

## 324 6 RESULTS

### 325 6.1 IS SELF-ADAPTIVE TRAINING LOSS A CALIBRATION LOSS?

326 To our understanding, the SAT method achieves the highest performance in the Selective Classification  
 327 problem with an unknown class. The first question to address is: *Is Self-Adaptive Training a  
 328 calibrated loss?* Since SAT adds an extra unknown class and modifies the loss function to alleviate  
 329 overfitting, it is reasonable to consider SAT loss as a potential calibration loss similar to Focal loss.  
 330 However, the role of SAT as a calibrator has not been explored in the literature. For this initial em-  
 331 pirical analysis, we set aside the Selective Classification problem and focus solely on the confidence  
 332 calibration problem.  
 333

334 A detailed analysis with graphs is presented in Appendix F. To sum up, first, we observed that the  
 335 accuracy and loss across epochs curves for Focal Loss exhibited similar trends, with minor overfit-  
 336 ting noted in the Food-101 dataset during the initial epochs. In contrast, the SAT loss demonstrated  
 337 different trends and did not consistently prevent overfitting. Regarding calibration metrics, the ECE  
 338 across epochs showed a consistent downward trend for Focal Loss, except for VGG-16 when applied  
 339 to the CIFAR-100 and Food-101 datasets, where an increase was observed in the initial epochs but  
 340 remained within an acceptable ECE range. For SAT loss and VGG-16, a rise was observed in ECE  
 341 after the 150 epochs in the first epochs decreasing after convergence, while the rise for ResNet-34  
 342 was less discernible. In the SAT end-to-end case, ECE values were notably high during the initial  
 343 epochs for both architectures. The MCE across epochs displayed similar trends for Focal Loss, but  
 344 distinct trends were observed for SAT. Importantly, SAT does not appear to be an effective calibra-  
 345 tion loss and may be detrimental when the goal is to train for a small number of epochs ( $\leq 100$ ), as  
 346 it outputs calibrated confidences only after a considerable amount of training epochs. Additionally,  
 347 we noted that the average confidence values of the *idk* class seem to be directly related to calibration,  
 348 reflecting similar trends as the ECE across epochs. This raised the question *Might the additional idk*  
 349 *class method be beneficial or detrimental in terms of calibration?* This observed behavior was the  
 350 main source of inspiration for incorporating predictions associated with the *idk* class into the novel  
 351 Socrates loss to calibrate during training.

352 Based on the empirical analysis the following claim can be made: Unlike Focal loss, which produces  
 353 very well-calibrated models and follows similar trends across all the datasets and architectures, SAT  
 354 loss exhibits certain tendencies that ultimately lead to the conclusion that it is not a loss that allows  
 355 learning calibrated models in all the epochs and scenarios, especially when aiming to train for a  
 356 small number of epochs or when dealing with complex datasets such as Food-101. Additionally,  
 357 when the loss is used *end-to-end*, the miscalibration in the first epochs is excessively large, and in  
 358 some cases (CIFAR-100 and Food-101 with VGG-16) it remains significantly large until the end of  
 359 training. When the loss is applied with *first-epochs* case, miscalibration begins to emerge. Therefore,  
 360 we can claim that SAT loss seems not to be a calibration loss.

### 361 6.2 SOCRATES LOSS AS A CALIBRATOR

362 Before addressing the topic of Selective Classification, a similar question asked in Section 6.1 needs  
 363 to be considered: *Is the novel Socrates loss a calibrated loss?* To investigate this, the same method-  
 364 ology of Section 6.1 is followed. Since SAT (*end-to-end* and *first-epochs*) has been empirically  
 365 shown not to be a suitable calibration loss, our novel method (*end-to-end*) is compared with the  
 366 CCL-SC method (*first-epochs* case as CCL-SC has two losses).  
 367

368 Due to space reasons, the curves for the SVHN and CIFAR-100 datasets, along with those for  
 369 CIFAR-10 and Food-101, are presented in Appendix G.  
 370

371 **Overfitting:** In contrast to the SAT method, both Socrates and CCL-SC effectively mitigate the  
 372 overfitting issue, improving generalization across all three datasets and both architectures. This  
 373 preliminary empirical analysis suggests that Socrates loss may be a prominent calibration loss. In  
 374 fact, the accuracies achieved with this novel loss outperform those obtained with the SAT loss,  
 375 showcasing a substantial improvement. When comparing Socrates with CCL-SC, the accuracies  
 376 are similar in most scenarios, except for the SVHN dataset with the VGG-16 architecture. Here,  
 377 Socrates achieves accuracies close to 100%, while CCL-SC reaches approximately 80%. Notably,  
 Socrates consistently exhibits a downward trend in output losses across all scenarios, whereas CCL-

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

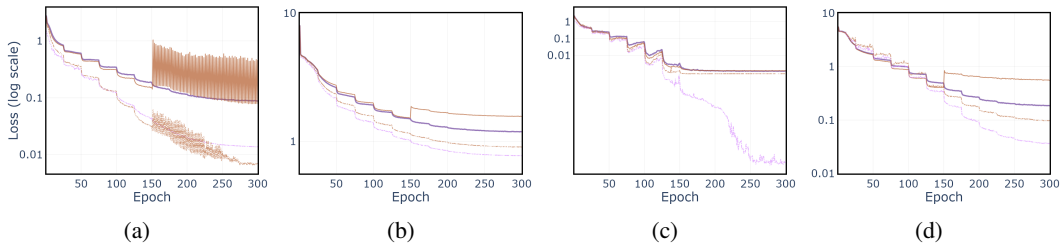


Figure 2: Loss curves of models trained on CIFAR-10 (a and c) and Food-101 (b and d) datasets using Socrates and CCL-SC methods with VGG-16 (a and b) and ResNet-34 (c and d) architectures.

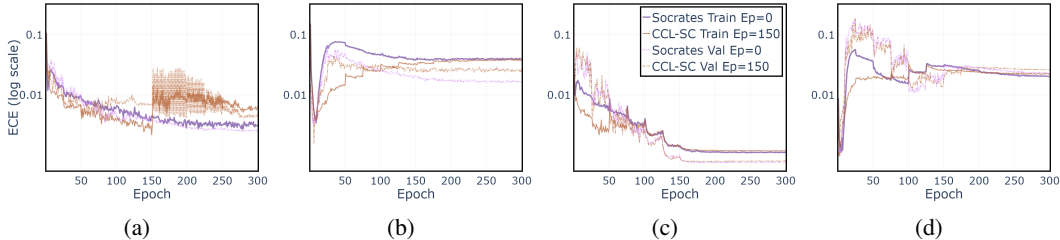


Figure 3: Evolution of the Expected Calibration Error (ECE) across epochs for models trained on CIFAR-10 (a and c) and Food-101 (b and d) datasets using Socrates and CCL-SC methods with VGG-16 (a and b) and ResNet-34 (c and d) architectures.

SC shows varied trends, including significant spikes and upward and downward trends depending on the architecture and dataset. Since both trainings use the same seeds and hyperparameters, except for the loss function, the observed spikes in CCL-SC suggest potential instability in training. The loss curves for CIFAR-10 and Food-101 are presented in Figure 2.

**Calibration Metrics:** The reliability diagrams with the ECE of the last epoch (Figure 1) do not provide enough information to draw calibration conclusions, instead, the ECE and MCE values along the epochs produce noticeable insights. In the first place, the ECE value along epochs is carried out. The values for CIFAR-10 and Food-101 are visualized in Figure 3.

Whereas SAT performs differently for each architecture and case, Socrates exhibits consistent trends across both architectures, showing a significant drop in ECE values after the initial epochs. Although Socrates shows an initial fluctuation in ECE values across epochs (which varies depending on the difficulty of the dataset), the ECE values across all epochs are consistently low, within a range below 10%, suggesting that the model is well-calibrated. Socrates achieves better ECE values than SAT across all epochs, datasets, and architectures.

When comparing Socrates with CCL-SC, both methods achieve similar ECE values. However, CCL-SC exhibits certain drawbacks. First, depending on the dataset and architecture, CCL-SC features spikes as equal to the loss across epochs. Second, while Socrates consistently shows an initial fluctuation followed by a decrease in ECE values, CCL-SC begins to miscalibrate the model once it reaches a lower ECE point. Although the ECE values remain within a small range, the upward trend indicates that CCL-SC could lead to miscalibrated models. Given that CCL-SC employs two losses (CCL loss for calibration and cross-entropy loss for classification), we argue that the calibration detriment could be attributed to the cross-entropy loss, which may miscalibrate the model by attempting to further reduce the loss after achieving the ideal confidence, thereby increasing the weight norm (Section 4.1.2). Therefore, having multiple losses could be detrimental if one of the losses is not specifically focused on calibration. This raises the question of why use several losses if a unified loss function can achieve the same goals, such as Socrates loss.

**Idk class:** Socrates reaffirms the claim made in Section 6.1: there is a link between the ECE values and the average of the confidences associated with the idk class. Moreover, addressing the question



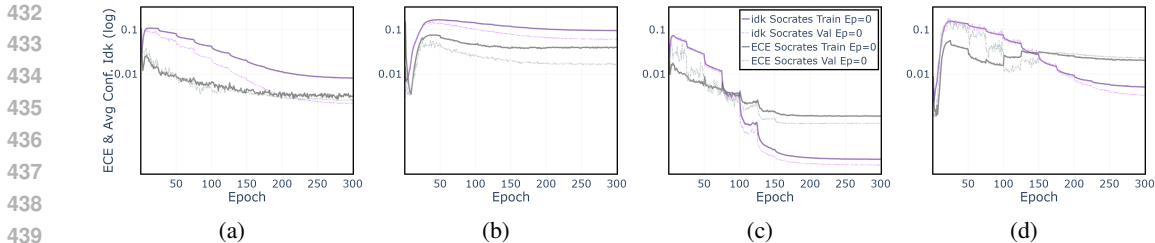


Figure 4: Curves depicting the average values of the *idk* class confidences across the epochs and ECE across epochs of models trained on CIFAR-10 (a and c) and Food-101 (b and d) datasets using Socrates method with VGG-16 (a and b) and ResNet-34 (c and d) architectures.

*Might the extra idk class approach be beneficial for calibration, or could it be detrimental?* we argue that incorporating the extra *idk* class and introducing  $\alpha_{dynamic}$  in the loss function offers a distinct advantage for calibrated selective classifiers. This mechanism, which is in Socrates loss, helps the model adjust penalization based on the confidence levels of its predictions. The curves showing the average values of the *idk* class confidences across epochs and the ECE across epochs for CIFAR-10 and Food-101 are presented in Figure 4.

**Socrates loss is a suitable loss to output calibrated models:** These findings underscore the effectiveness of Socrates as an end-to-end calibration method for training models, particularly when only a small number of epochs (in contrast to SAT, which is not suitable) are required to train trustworthy outputs in terms of confidence calibration and when an unknown class is considered. The ability of Socrates to function without having multiple losses allows for a unified optimization goal, simultaneously addressing both classification and calibration in an ad-hoc manner.

### 6.2.1 SOCRATES LOSS AS A SELECTIVE CLASSIFIER

This paper focuses on calibrating selective classifiers, aiming to produce ad-hoc calibrated selective classifiers suitable for deployment in real-world critical environments. While improving Selective Classification error rates was not the primary goal of our study, which was more focused on enhancing calibration, Socrates demonstrates improvements over SAT on challenging datasets such as Food-101. In comparison with CCL-SC, both methods achieve comparable performance. There are instances where Socrates outperforms, as in SVHN with VGG-16, where Socrates achieves an error rate close to 3% compared to around 18% for CCL-SC and SAT for the *first-epochs* case. The most relevant results are presented in Table 2, and for space reasons in Appendix E. The risk-coverage curves provide a clear demonstration of the strength of the Socrates method compared to CCL-SC. These curves reveal that Socrates consistently achieves similar or better values than CCL-SC. The detailed Risk-Coverage curves can be found in Figure 5 and Appendix H, where a notable improvement is observed, particularly for the CIFAR-10 and SVHN datasets.

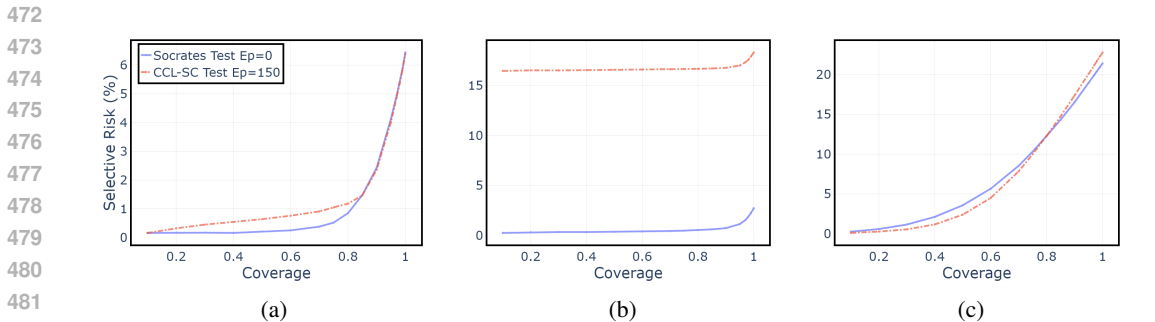


Figure 5: Risk-Coverage curves of models trained on CIFAR-10 (a), SVHN (b) and Food-101 (c) datasets using Socrates (*end-to-end* case) and CCL-SC (*first-epochs* case) methods with VGG-16 (a and b) and ResNet-34 (c) architectures.

Table 2: Selective Classification error rates % on CIFAR-10, SVHN and Food-101 datasets for various coverage rates %, reported with mean and standard deviation. Underline indicate the overall best performance, while bold highlight the best performance in each case.

Dataset	Coverage	<i>end-to-end</i> case		<i>first-epochs</i> case		
		Socrates (ours)	SAT	Socrates + Focal	CCL-SC + c.e	SAT + c.e
CIFAR-10	100	<b>6.44 ± 0.18</b>	7.08 ± 1.07	6.67 ± 0.19	<u>6.38 ± 0.14</u>	6.87 ± 1.08
	95	<b>4.14 ± 0.12</b>	4.78 ± 0.98	4.45 ± 0.14	<u>4.02 ± 0.14</u>	4.58 ± 1.12
	90	<b>2.43 ± 0.09</b>	3.01 ± 0.88	2.76 ± 0.13	<u>2.36 ± 0.13</u>	2.92 ± 1.01
	85	<b>1.48 ± 0.11</b>	1.82 ± 0.65	1.64 ± 0.20	<u>1.47 ± 0.16</u>	1.75 ± 0.74
	80	<b>0.85 ± 0.03</b>	1.12 ± 0.51	<b>1.05 ± 0.11</b>	1.18 ± 0.25	1.05 ± 0.46
	75	<b>0.52 ± 0.03</b>	0.67 ± 0.32	<i>0.68 ± 0.07</i>	1.05 ± 0.19	<b>0.61 ± 0.27</b>
	70	<b>0.38 ± 0.04</b>	0.43 ± 0.24	<i>0.51 ± 0.05</i>	0.91 ± 0.11	<b>0.42 ± 0.20</b>
SVHN	100	2.72 ± 0.07	<b>2.65 ± 0.04</b>	<b>2.80 ± 0.03</b>	18.29 ± 34.73	18.22 ± 34.77
	95	1.15 ± 0.04	<b>1.04 ± 0.02</b>	<b>1.20 ± 0.08</b>	16.99 ± 35.46	16.89 ± 35.51
	90	0.74 ± 0.05	<b>0.61 ± 0.05</b>	<b>0.80 ± 0.05</b>	16.76 ± 35.58	16.57 ± 35.69
	85	0.62 ± 0.02	<b>0.45 ± 0.04</b>	<b>0.62 ± 0.05</b>	16.70 ± 35.62	16.44 ± 35.76
	80	0.55 ± 0.03	<b>0.38 ± 0.02</b>	<b>0.54 ± 0.05</b>	16.66 ± 35.64	16.39 ± 35.79
	75	0.49 ± 0.05	<b>0.33 ± 0.02</b>	<b>0.51 ± 0.03</b>	16.64 ± 35.65	16.35 ± 35.81
	70	0.45 ± 0.04	<b>0.30 ± 0.01</b>	<b>0.48 ± 0.02</b>	16.62 ± 35.66	16.33 ± 35.82
Food-101	100	<b>21.40 ± 0.79</b>	100 ± 0.0	32.33 ± 22.32	22.77 ± 0.90	<b>22.08 ± 0.75</b>
	95	<b>18.95 ± 0.80</b>	100 ± 0.0	30.20 ± 23.10	20.09 ± 0.92	<b>20.02 ± 0.74</b>
	90	<b>16.54 ± 0.75</b>	100 ± 0.0	28.23 ± 23.92	<b>17.39 ± 0.91</b>	17.97 ± 0.74
	85	<b>14.32 ± 0.74</b>	100 ± 0.0	26.37 ± 23.92	<b>14.75 ± 0.92</b>	15.99 ± 0.72
	80	<b>12.30 ± 0.78</b>	100 ± 0.0	24.60 ± 25.11	<b>12.30 ± 0.94</b>	14.08 ± 0.67
	75	<b>10.32 ± 0.68</b>	100 ± 0.0	22.94 ± 25.57	<b>10.00 ± 0.81</b>	12.20 ± 0.64
	70	<b>8.54 ± 0.62</b>	100 ± 0.0	21.49 ± 25.97	<b>7.85 ± 0.70</b>	10.37 ± 0.60

## 7 CONCLUSIONS AND LIMITATIONS

In this paper, we first empirically investigated the calibration capacity of SAT loss as a calibration mechanism, finding that it does not produce well-calibrated models. This deficiency is particularly detrimental for models that require only a small number of epochs or when working with hard-to-classify datasets. Additionally, we found that SAT does not consistently mitigate overfitting across all cases. Through this empirical study, we identified a relationship between the extra unknown class and calibration, which inspired the development of our proposed loss function. To address the need for ad-hoc easy-to-implement calibrated selective classifiers with an unknown class, we proposed a new method that incorporates an extra unknown class and introduces a novel loss, *Socrates*, with a unified optimization goal (classification and calibration). We theoretically and empirically analyzed this loss, demonstrating that it is an optimal calibration method without the previously enumerated drawbacks of SAT and CCL-SC. This new loss not only ensures strong calibration throughout all training epochs (making it suitable for models trained with fewer epochs), but also produces selective classifiers that achieve similar Selective Classification error rates to SAT and CCL-SC, while notably outperforming SAT on hard-to-classify datasets such as CIFAR-100 and Food-101 and CCL-SC on datasets such as SVHN for VGG-16.

We encourage the research community to further evaluate the Socrates method across a broader spectrum of architectures and datasets. It is notable that this method has not been compared to post-hoc methods. We argue that the strength of this end-to-end method comes from producing compact models that do not require post-processing and additional data. Leveraging all available data during training can be particularly advantageous when data is limited. Future research should incorporate and evaluate additional reliability aspects to develop a more comprehensive reliability framework (e.g., distribution shifts, noise, out-of-distribution, etc.). The lack of metrics that integrate reliability concepts is a pressing need. For example, a model may often be well-calibrated but exhibit low accuracy. Additionally, there is a need for metrics that summarize calibration performance across epochs. For example, a new ECE-epochs metric could indicate whether calibration has improved or deteriorated at any given point.

## REFERENCES

- Abiodun Ayodeji, Muritala Alade Amidu, Samuel Abiodun Olatubosun, Yacine Addad, and Hafiz Ahmed. Deep learning for safety assessment of nuclear power reactors: Reliability, explainability, and research opportunities. *Progress in Nuclear Energy*, 151:104339, 2022. ISSN 0149-1970. doi: <https://doi.org/10.1016/j.pnucene.2022.104339>. URL <https://www.sciencedirect.com/science/article/pii/S0149197022002141>.
- Abhijit Bendale and Terrance Boulton. Towards open set deep networks. In *Computer Vision and Pattern Recognition (CVPR), 2016 IEEE Conference on*. IEEE, 2016.
- Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101 – mining discriminative components with random forests. In David Fleet, Tomas Pajdla, Bernt Schiele, and Tinne Tuytelaars (eds.), *Computer Vision – ECCV 2014*, pp. 446–461, Cham, 2014. Springer International Publishing. ISBN 978-3-319-10599-4.
- Corinna Cortes, Giulia DeSalvo, and Mehryar Mohri. Learning with rejection. In Ronald Ortner, Hans Ulrich Simon, and Sandra Zilles (eds.), *Algorithmic Learning Theory*, pp. 67–82, Cham, 2016. Springer International Publishing. ISBN 978-3-319-46379-7.
- Ran El-Yaniv and Yair Wiener. On the foundations of noise-free selective classification. *Journal of Machine Learning Research*, 11(53):1605–1641, 2010. URL <http://jmlr.org/papers/v11/el-yaniv10a.html>.
- Leo Feng, Mohamed Osama Ahmed, Hossein Hajimirsadeghi, and Amir H. Abdi. Towards better selective classification. In *The Eleventh International Conference on Learning Representations*, 2023. URL [https://openreview.net/forum?id=5gDz\\_yTcst](https://openreview.net/forum?id=5gDz_yTcst).
- Telmo Silva Filho, Hao Song, Miquel Perello-Nieto, Raul Santos-Rodriguez, Meelis Kull, and Peter Flach. Classifier Calibration: A survey on how to assess and improve predicted class probabilities. *Machine Learning*, 112(9):3211–3260, September 2023. ISSN 0885-6125, 1573-0565. doi: 10.1007/s10994-023-06336-7. URL <https://doi.org/10.1007/s10994-023-06336-7>.
- Adam Fisch, Tommi S. Jaakkola, and Regina Barzilay. Calibrated selective classification. *Transactions on Machine Learning Research*, 2022. ISSN 2835-8856. URL <https://openreview.net/forum?id=zFhNBs8GaV>.
- Ido Galil, Mohammed Dabbah, and Ran El-Yaniv. What can we learn from the selective prediction and uncertainty estimation performance of 523 imagenet classifiers? In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=p66AzKi6Xim>.
- Yonatan Geifman and Ran El-Yaniv. Selective classification for deep neural networks. In *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS’17*, pp. 4885–4894, Red Hook, NY, USA, 2017. Curran Associates Inc. ISBN 9781510860964.
- Yonatan Geifman and Ran El-Yaniv. Selectivenet: A deep neural network with an integrated reject option. In *International Conference on Machine Learning*, 2019. URL <https://api.semanticscholar.org/CorpusID:59316904>.
- Elakktat D. Gireesh and Varadaraj P. Gurupur. Information entropy measures for evaluation of reliability of deep neural network results. *Entropy*, 25(4), 2023. ISSN 1099-4300. doi: 10.3390/e25040573. URL <https://www.mdpi.com/1099-4300/25/4/573>.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On Calibration of Modern Neural Networks. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70, ICML’17*, pp. 1321–1330. JMLR.org, 2017.
- Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *Proceedings of the International Conference on Learning Representations*, 2019.

- 594 José Hernández-Orallo, Peter Flach, and Cèsar Ferri. A unified view of performance met-  
595 rics: Translating threshold choice into expected classification loss. *Journal of Machine*  
596 *Learning Research*, 13(91):2813–2869, 2012. URL [http://jmlr.org/papers/v13/](http://jmlr.org/papers/v13/hernandez-orallo12a.html)  
597 [hernandez-orallo12a.html](http://jmlr.org/papers/v13/hernandez-orallo12a.html).
- 598
- 599 Lang Huang, Chao Zhang, and Hongyang Zhang. Self-Adaptive Training: beyond Empirical Risk  
600 Minimization. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.),  
601 *Advances in Neural Information Processing Systems*, volume 33, pp. 19365–19376. Curran  
602 Associates, Inc., 2020. URL [https://proceedings.neurips.cc/paper\\_files/](https://proceedings.neurips.cc/paper_files/paper/2020/file/e0ab531ec312161511493b002f9be2ee-Paper.pdf)  
603 [paper/2020/file/e0ab531ec312161511493b002f9be2ee-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/e0ab531ec312161511493b002f9be2ee-Paper.pdf).
- 604 Alex Krizhevsky. Learning multiple layers of features from tiny images. *Technical Report Uni-*  
605 *versity of Toronto*, pp. 32–33, 2009. URL [https://www.cs.toronto.edu/~kriz/](https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf)  
606 [learning-features-2009-TR.pdf](https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf).
- 607
- 608 Yann LeCun, Bernhard E. Boser, John S. Denker, Donnie Henderson, Richard E. Howard,  
609 Wayne E. Hubbard, and Lawrence D. Jackel. Handwritten digit recognition with a back-  
610 propagation network. In *Neural Information Processing Systems*, 1989. URL [https://api.](https://api.semanticscholar.org/CorpusID:2542741)  
611 [semanticscholar.org/CorpusID:2542741](https://api.semanticscholar.org/CorpusID:2542741).
- 612 Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollar. Focal Loss for Dense Ob-  
613 ject Detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(2):318–327,  
614 February 2020. ISSN 0162-8828, 2160-9292, 1939-3539. doi: 10.1109/TPAMI.2018.2858826.  
615 URL <https://ieeexplore.ieee.org/document/8417976/>.
- 616
- 617 Bingyuan Liu, Jérôme Rony, Adrian Galdran, Jose Dolz, and Ismail Ben Ayed. Class Adaptive  
618 Network Calibration. In *CVPR*, April 2023. URL <http://arxiv.org/abs/2211.15088>.
- 619
- 620 Ziyin Liu, Zhikang Wang, Paul Pu Liang, Russ R Salakhutdinov, Louis-Philippe Morency,  
621 and Masahito Ueda. Deep gamblers: Learning to abstain with portfolio theory. In  
622 H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett (eds.),  
623 *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.,  
624 2019. URL [https://proceedings.neurips.cc/paper\\_files/paper/2019/](https://proceedings.neurips.cc/paper_files/paper/2019/file/0c4b1eeb45c90b52bfb9d07943d855ab-Paper.pdf)  
625 [file/0c4b1eeb45c90b52bfb9d07943d855ab-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/0c4b1eeb45c90b52bfb9d07943d855ab-Paper.pdf).
- 626 Ben McEwen, Richard Green, Stefanie Gutschmidt, and Grant Ryan. Predictive state estima-  
627 tion of invasive predators using low resolution thermal cameras. In *2021 36th International*  
628 *Conference on Image and Vision Computing New Zealand (IVCNZ)*, pp. 1–6, 2021. doi:  
629 10.1109/IVCNZ54163.2021.9653201.
- 630
- 631 Jooyoung Moon, Jihyo Kim, Younghak Shin, and Sangheum Hwang. Confidence-aware learning for  
632 deep neural networks. In *Proceedings of the 37th International Conference on Machine Learning*,  
633 ICML’20. JMLR.org, 2020.
- 634 Jishnu Mukhoti, Viveka Kulharia, Amartya Sanyal, Stuart Golodetz, Philip H. S. Torr, and Puneet K.  
635 Dokania. Calibrating deep neural networks using focal loss. In *Proceedings of the 34th Inter-*  
636 *national Conference on Neural Information Processing Systems*, NIPS’20, Red Hook, NY, USA,  
637 2020. Curran Associates Inc. ISBN 9781713829546.
- 638
- 639 Mahdi Pakdaman Naeini, Gregory F. Cooper, and Milos Hauskrecht. Obtaining well calibrated  
640 probabilities using Bayesian Binning. *Proceedings of the National Conference on Artificial Intel-*  
641 *ligence*, 4:2901–2907, 2015. ISSN 2159-5399. doi: 10.1609/aaai.v29i1.9602.
- 642 Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y. Ng. Reading  
643 digits in natural images with unsupervised feature learning. In *NIPS Workshop on Deep Learning*  
644 *and Unsupervised Feature Learning 2011*, 2011. URL [http://ufldl.stanford.edu/](http://ufldl.stanford.edu/housenumbers/nips2011_housenumbers.pdf)  
645 [housenumbers/nips2011\\_housenumbers.pdf](http://ufldl.stanford.edu/housenumbers/nips2011_housenumbers.pdf).
- 646
- 647 Khanh Nguyen and Brendan O’Connor. Posterior calibration and exploratory analysis for natural  
language processing models. *Proceedings of EMNLP*, 2015. ISSN 2331-8422.

- 648 Alexandru Niculescu-Mizil and Rich Caruana. Predicting good probabilities with supervised learn-  
649 ing. In *ACM International Conference Proceeding Series; Vol. 119: Proceedings of the 22nd in-*  
650 *ternational conference on Machine learning; 07-11 Aug. 2005*, pp. 625–632. ACM, 2005. ISBN  
651 1595931805.
- 652 Yiwen Liao Patrick Schlachter and Bin Yang. Open-set recognition using intra-class splitting. In  
653 *2019 IEEE European Signal Processing Conference (EUSIPCO)*, September 2019.
- 654 Gabriel Pereyra, George Tucker, Jan Chorowski, Łukasz Kaiser, and Geoffrey Hinton. Regularizing  
655 neural networks by penalizing confident output distributions, 2017. URL [https://arxiv.](https://arxiv.org/abs/1701.06548)  
656 [org/abs/1701.06548](https://arxiv.org/abs/1701.06548).
- 657 J. Platt. Probabilistic outputs for support vector machines and comparison to regularized likelihood  
658 methods. In *Advances in Large Margin Classifiers*, 2000.
- 659 Dustin Tran, Jeremiah Liu, Michael W. Dusenberry, Du Phan, Mark Collier, Jie Ren, Kehang Han,  
660 Zi Wang, Zeldia Mariet, Huiyi Hu, Neil Band, Tim G. J. Rudner, Karan Singhal, Zachary Nado,  
661 Joost van Amersfoort, Andreas Kirsch, Rodolphe Jenatton, Nithum Thain, Honglin Yuan, Kelly  
662 Buchanan, Kevin Murphy, D. Sculley, Yarin Gal, Zoubin Ghahramani, Jasper Snoek, and Bal-  
663 aji Lakshminarayanan. Plex: Towards Reliability using Pretrained Large Model Extensions. In  
664 *Proceedings of the International conference on Machine Learning - ICML '22. Pre-training Work-*  
665 *shop*, July 2022. URL <http://arxiv.org/abs/2207.07411>.
- 666 Yu-Chang Wu, Shen-Huan Lyu, Haopu Shang, Xiangyu Wang, and Chao Qian. Confidence-aware  
667 contrastive learning for selective classification. In *International Conference on Machine Learning*,  
668 2024.
- 669 Xu-Yao Zhang, Guo-Sen Xie, Xiuli Li, Tao Mei, and Cheng-Lin Liu. A Survey on Learn-  
670 ing to Reject. *Proceedings of the IEEE*, 111(2):185–215, February 2023. ISSN 0018-9219,  
671 1558-2256. doi: 10.1109/JPROC.2023.3238024. URL [https://ieeexplore.ieee.org/](https://ieeexplore.ieee.org/document/10028760/)  
672 [document/10028760/](https://ieeexplore.ieee.org/document/10028760/).
- 673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

APPENDIX: TRAINING RELIABLE MODELS:  
HAVING THE CONFIDENCE TO SAY “I DON’T KNOW”

A PSEUDOCODE FOR THE SOCRATES METHOD

---

**Algorithm 1** Training with Socrates loss

---

**Require:** Data  $\{(x_i, y_i)\}_{i=1}^n$ , initial targets  $\{t_i\}_{i=1}^n = \{y_i\}_{i=1}^n$ , initial model  $f$ , batch size  $bs$ , and hyper-parameters: momentum term  $\alpha_{\text{momentum}}$ , modularity factor  $\gamma$ , and initial epochs  $E_s$ .

```

1: repeat
2:   for e = 1 to maximum_epochs do
3:     for each mini-batch data  $\{(x_i, y_i)\}_{bs}$  in the current epoch e do
4:       for i = 1 to bs (in parallel) do
5:          $\hat{p}_i = \text{softmax}(f(x_i))$ 
6:          $\alpha_{\text{dynamic}} \leftarrow (\max_{y_i \neq y_{gt}} \hat{p}_{i, y_i}) - \hat{p}_{i, \text{idk}}$ 
7:         if  $e \geq E_s$  then
8:            $t_i \leftarrow \alpha_{\text{momentum}} \times t_i + (1 - \alpha_{\text{momentum}}) \times \hat{p}_{i, y_i}$ 
9:         end if
10:         $\mathcal{L}(f) = -\frac{1}{n} \sum_{i=1}^n (1 - \hat{p}_{i, y_i})^\gamma [t_{i, y_i} \log \hat{p}_{i, y_i} + \alpha_{\text{dynamic}} (1 - t_{i, y_i}) \log \hat{p}_{i, \text{idk}}]$ 
11:        Update the weights of  $f$  using an optimizer based on  $\mathcal{L}(f)$ 
12:      end for
13:    end for
14:  end for
15: until end of training

```

---

Although our method can be used with other losses due to the flexibility of the initial epochs variable, our primary goal is to design an end-to-end loss. Therefore, we set  $E_s = 0$  in our main results.

B MATHEMATICAL EXAMPLE OF THE SOCRATES METHOD

To illustrate how Socrates loss operates, consider a selective classifier with  $E_s = 0$ ,  $\gamma = 2$ , and  $\alpha_{\text{momentum}} = 0.9$ , which can output one of three classes: predator, non-predator, or idk. We will examine the following three scenarios:

1. An image of a cat with a ground truth label of predator. The classifier outputs the confidences  $[0.9, 0.05, 0.05]$  at epoch 30, and  $[0.9, 0.02, 0.08]$  at epoch 31. Since the previous prediction (epoch 30) had high confidence, the  $t_i = 0.9$ . For epoch 31, as  $\max_{y_i \neq y_{gt}} \hat{p}_{i, y_i}$  corresponds to the idk class,  $\alpha_{\text{dynamic}} = 0$ . Therefore, only the first part of the loss function is relevant, giving more penalty to hard-to-classify instances. The loss at epoch 31 is  $\mathcal{L} = 0.0009$ .
2. An image of a fake pink cat with a ground truth label of predator. The classifier outputs the confidences  $[0.5, 0.25, 0.25]$  at epoch 30, and  $[0.5, 0.3, 0.2]$  at epoch 31. Since the previous prediction (epoch 30) lacked high confidence, both parts of the equation are relevant,  $t_i = 0.5$ . In this case,  $\max_{y_i \neq y_{gt}} \hat{p}_{i, y_i}$  is the non-predator class, then  $\alpha_{\text{dynamic}} = 0.1$ ; the model is unaware of its lack of knowledge. The loss at epoch 31 is  $\mathcal{L} = 0.11$ .
3. An image of a fake pink cat with a ground truth label of predator. The classifier outputs the confidences  $[0.5, 0.25, 0.25]$  at epoch 30, and  $[0.5, 0.2, 0.3]$  at epoch 31. As the previous prediction (epoch 30) lacked high confidence, both parts of the equation take relevance,  $t_i = 0.5$ . In this case,  $\max_{y_i \neq y_{gt}} \hat{p}_{i, y_i}$  is the idk class, then  $\alpha_{\text{dynamic}} = 0$ ; the model is aware of its lack of knowledge. The loss at epoch 31 is  $\mathcal{L} = 0.088$ .

## C THEORETICAL PROOFS

### C.1 SOCRATES LOSS FORMS A REGULARIZED UPPER BOUND IN THE KULLBACK-LEIBLER DIVERGENCE

**Theorem:** Socrates loss minimizes (creates an upper bound for) the Kullback-Leibler (KL) divergence while regularizing by increasing the entropy of the predicted distribution and leveraging the predictions associated with the unknown class. The regularization parameters are  $\gamma$ ,  $\alpha_{\text{dynamic}}$ , and  $\Delta_{\text{reg}}$ ; where  $\Delta_{\text{reg}} = (1 - t_y)[\gamma \hat{p}_y \log \hat{p}_{idk} - \log \hat{p}_{idk}]$ . Therefore:

$$\mathcal{L}(f) \geq D_{\text{KL}}(q||\hat{p}) - \gamma \mathbb{H}[\hat{p}] + \alpha_{\text{dynamic}} \Delta_{\text{reg}}; \quad (4)$$

**Proof:** Let the KL divergence be the divergence between the ground truth distribution  $q$  and the predicted distribution  $\hat{p}$ , and  $\mathbb{H}[q]$  be the entropy of the ground truth distribution defined as  $\mathbb{H}[q] = -\sum_j q_j \log(q_j)$ . Therefore, for a multiclass problem, the KL divergence can be expressed as:

$$\begin{aligned} D_{\text{KL}}(q||\hat{p}) &= \sum_j q_j \log\left(\frac{q_j}{\hat{p}_j}\right) = \\ &= \sum_j q_j \log(q_j) - \sum_j q_j \log(\hat{p}_j); \Rightarrow \\ &\Rightarrow D_{\text{KL}}(q||\hat{p}) = -\mathbb{H}[q] + \mathcal{L}_{\text{c.e.}}(f); \end{aligned} \quad (5)$$

where  $\mathcal{L}_{\text{c.e.}}(f)$  is the cross-entropy loss, which forms an upper bound in the KL divergence:

$$\begin{aligned} \mathcal{L}_{\text{c.e.}}(f) &= D_{\text{KL}}(q||\hat{p}) + \mathbb{H}[q]; \Rightarrow \\ &\Rightarrow \mathcal{L}_{\text{c.e.}}(f) \geq D_{\text{KL}}(q||\hat{p}); \end{aligned} \quad (6)$$

To simplify, we consider the case of the first selected epochs where  $t_i \leftarrow y_i = 1$ . Let  $t_i \in q$ , be the target distribution. If we take only one instance of  $m$  number of instances, i.e.  $m = 1$ , the loss function can be written as:

$$\mathcal{L}(f) = -[t_y(1 - \hat{p}_y)^\gamma \log \hat{p}_y + \alpha_{\text{dynamic}}(1 - t_y)(1 - \hat{p}_y)^\gamma \log \hat{p}_{idk}], \quad (7)$$

where the subscript  $y$  denotes the values associated with the ground truth class.

Using Bernoulli's inequality, which states that  $(1 - x)^\alpha \geq 1 - \alpha x$ , if  $0 \leq x \leq 1$  and  $\alpha \geq 0$ , as  $\forall \gamma \geq 1$  and the  $\hat{p}_y \in [0, 1]$ , then we get:

$$\begin{aligned} \mathcal{L}(f) &= -(1 - \hat{p}_y)^\gamma [t_y \log \hat{p}_y + \alpha_{\text{dynamic}}(1 - t_y) \log \hat{p}_{idk}] \\ &\geq -(1 - \gamma \hat{p}_y) [t_y \log \hat{p}_y + \alpha_{\text{dynamic}}(1 - t_y) \log \hat{p}_{idk}] \\ &= \gamma \hat{p}_y t_y \log \hat{p}_y - t_y \log \hat{p}_y + \gamma \hat{p}_y \alpha_{\text{dynamic}}(1 - t_y) \log \hat{p}_{idk} - \alpha_{\text{dynamic}}(1 - t_y) \log \hat{p}_{idk} \\ &= -\gamma \mathbb{H}[\hat{p}] + \mathcal{L}_{\text{c.e.}}(f) + \alpha_{\text{dynamic}} \Delta_{\text{reg}} \\ &= -\gamma \mathbb{H}[\hat{p}] + D_{\text{KL}}(q||\hat{p}) + \mathbb{H}[q] + \alpha_{\text{dynamic}} \Delta_{\text{reg}}; \\ &\quad \text{where } \Delta_{\text{reg}} = (1 - t_y)[\gamma \hat{p}_y \log \hat{p}_{idk} - \log \hat{p}_{idk}]; \end{aligned} \quad (8)$$

Therefore:

$$\mathcal{L}(f) \geq D_{\text{KL}}(q||\hat{p}) + \mathbb{H}[q] - \gamma \mathbb{H}[\hat{p}] + \alpha_{\text{dynamic}} \Delta_{\text{reg}}; \quad (9)$$

where  $\mathbb{H}[q]$  is a constant.

Thus, this new loss improves calibration by minimizing the KL divergence, maximizing the entropy depending on the weight of  $\gamma$  (which smooths the learned distributions), and adding an extra regularization term (which might help to avoid overfitting) which maximises the uncertainty when the prediction is incorrect.

### C.2 SOCRATES LOSS REGULARIZES THE WEIGHTS OF THE NETWORK

Let  $\mathcal{L}_{\text{c.e.}}(f)$  be cross-entropy loss, and  $\mathcal{L}(f)$  be Socrates loss. The gradients of the neural network trained with  $\mathcal{L}(f)$  are smaller than the ones trained with  $\mathcal{L}_{\text{c.e.}}(f)$  when perfect confidence is reached and the model could start overfitting and subsequently be miscalibrated, i.e.,

$$\left\| \frac{\partial \mathcal{L}(f)}{\partial w} \right\| \leq \left\| \frac{\partial \mathcal{L}_{\text{c.e.}}(f)}{\partial w} \right\|. \quad (10)$$

This behaviour shows that Socrates loss acts as a regularizer when the model is sufficiently confident, avoiding miscalibration and overfitting.

**Proof:** To simplify, we consider the case of the first selected epochs where  $t_i \leftarrow y_i = 1$ . If we take only one instance from  $m$  instances, i.e.  $m = 1$ , the Socrates loss function can be written as:

$$\mathcal{L}(f) = -[t_y(1 - \hat{p}_y)^\gamma \log \hat{p}_y + \alpha_{\text{dynamic}}(1 - t_y)(1 - \hat{p}_y)^\gamma \log \hat{p}_{idk}]. \quad (11)$$

The gradient with respect to the parameters of the last linear layer can be decomposed with the chain rule:

$$\begin{aligned} \frac{\partial \mathcal{L}(f)}{\partial w} &= \frac{\partial \mathcal{L}(f)}{\partial \hat{p}_y} \frac{\partial \hat{p}_y}{\partial z} \frac{\partial z}{\partial w}. \\ \text{where } \frac{\partial \mathcal{L}(f)}{\partial \hat{p}_y} &= \gamma(1 - \hat{p}_y)^{\gamma-1} t_y \log \hat{p}_y - (1 - \hat{p}_y)^\gamma \frac{t_y}{\hat{p}_y} + \\ &+ \gamma(1 - \hat{p}_y)^{\gamma-1} \alpha_{\text{dynamic}}(1 - t_y) \log \hat{p}_{idk} - (1 - \hat{p}_y)^\gamma \alpha_{\text{dynamic}}(1 - t_y) \frac{1}{\hat{p}_{idk}}. \end{aligned} \quad (12)$$

On the other hand, cross-entropy loss can be written as:

$$\mathcal{L}_{c.e.}(f) = -t_y \log \hat{p}_y. \quad (13)$$

Where the gradient using the chain rule is:

$$\begin{aligned} \frac{\partial \mathcal{L}_{c.e.}(f)}{\partial w} &= \frac{\partial \mathcal{L}_{c.e.}(f)}{\partial \hat{p}_y} \frac{\partial \hat{p}_y}{\partial z} \frac{\partial z}{\partial w}. \\ \text{where } \frac{\partial \mathcal{L}_{c.e.}(f)}{\partial \hat{p}_y} &= -\frac{t_y}{\hat{p}_y} \end{aligned} \quad (14)$$

Then, we can observe that the gradient of cross-entropy is a component of the gradient of Socrates:

$$\begin{aligned} \frac{\partial \mathcal{L}(f)}{\partial \hat{p}_y} &= \frac{\partial \mathcal{L}_{c.e.}(f)}{\partial \hat{p}_y} [(1 - \hat{p}_y)^\gamma - \gamma \hat{p}_y (1 - \hat{p}_y)^{\gamma-1} \log \hat{p}_y] + \\ &+ \gamma(1 - \hat{p}_y)^{\gamma-1} \alpha_{\text{dynamic}}(1 - t_y) \log \hat{p}_{idk} - (1 - \hat{p}_y)^\gamma \alpha_{\text{dynamic}}(1 - t_y) \frac{1}{\hat{p}_{idk}}. \end{aligned} \quad (15)$$

If  $g(\hat{p}_y, \gamma) = (1 - \hat{p}_y)^\gamma - \gamma \hat{p}_y (1 - \hat{p}_y)^{\gamma-1} \log \hat{p}_y$  is a regularizer of the cross-entropy; and  $r(t_y, \alpha_{\text{dynamic}}, \hat{p}_y, \hat{p}_{idk}) = \gamma(1 - \hat{p}_y)^{\gamma-1} \alpha_{\text{dynamic}}(1 - t_y) \log \hat{p}_{idk} - (1 - \hat{p}_y)^\gamma \alpha_{\text{dynamic}}(1 - t_y) \frac{1}{\hat{p}_{idk}}$  is highly affected by the idk class, which adds a small penalty  $r(t_y, \alpha_{\text{dynamic}}, \hat{p}_y, \hat{p}_{idk}) \in [0, 1]$ , then:

$$\frac{\partial \mathcal{L}(f)}{\partial \hat{p}_y} = \frac{\partial \mathcal{L}_{c.e.}(f)}{\partial \hat{p}_y} g(\hat{p}_y, \gamma) + r(t_y, \alpha_{\text{dynamic}}, \hat{p}_y, \hat{p}_{idk}). \quad (16)$$

When confidence is high, and the model could start being overfitted and miscalibrated, the value of  $g(\hat{p}_y, \gamma) \in [0, 1]$ . In that case:

$$\left\| \frac{\partial \mathcal{L}(f)}{\partial \hat{p}_y} \right\| \leq \left\| \frac{\partial \mathcal{L}_{c.e.}(f)}{\partial \hat{p}_y} \right\| \implies \left\| \frac{\partial \mathcal{L}(f)}{\partial w} \right\| \leq \left\| \frac{\partial \mathcal{L}_{c.e.}(f)}{\partial w} \right\| \quad (17)$$

This demonstrates that the gradients of a model associated with the Socrates loss are smaller than those associated with the cross-entropy loss when perfect confidence is reached. Therefore, the Socrates loss acts as a regularizer with a penalty associated with the unknown knowledge of the classifier, avoiding overfitting, and subsequently miscalibration.

## D MODEL REPRODUCIBILITY

### D.1 COMPUTE

The experiments were conducted on a shared supercomputer (Nvidia A100 80Gb SXM4 GPU). We consider it inequitable to provide specific time allocations for each method due to the nature of a shared supercomputer, where training durations vary based on resource availability. To ensure fair results, five different seeds were employed for each method, case, dataset, and architecture. The list of seeds to replicate results can be found in Table 3.



With certain methods and seeds, the training failed to achieve high accuracies, remaining stuck from the start at levels close to 10 and 20%. For SVHN, Focal loss could not train (converge) with VGG-16 when using seed 403. Additionally for SVHN, SAT and CCL-SC failed to train with VGG-16 using seed 303, and CCL-SC failed with VGG-16 for seeds 402, 403, 404, 405, and 409. For Food-101, SAT was unable to train with VGG-16 and ResNet-34 for any  $E(s) = 0$  seed. In contrast, Socrates successfully trained under all conditions.

Table 3: Seeds for results replication

Dataset	Es	Seeds VGG-16	Seeds ResNet-34
CIFAR-10 and CIFAR-10C	150	301, 302, 303, 304, 309	305, 306, 307, 308, 309
	0	401, 402, 403, 404, 409	405, 406, 407, 408, 409
CIFAR-100 and CIFAR-100C	150	301, 302, 303, 304, 309	305, 306, 307, 308, 309
	0	401, 402, 403, 404, 409	405, 406, 407, 408, 409
SVHN	150	301, 302, 303, 304, 309	305, 306, 307, 308, 309
	0	401, 402, 403, 404, 409	405, 406, 407, 408, 409
Food-101	150	301, 312, 313, 314, 319	311, 312, 313, 314, 319
	0	401, 412, 413, 414, 419	411, 412, 413, 414, 419

## D.2 HYPERPARAMETERS

To conduct the experiments, we adapted the publicly available official implementation of Self-Adaptative Training, which was adapted from DeepGamblers (Liu et al., 2019). The hyperparameter values do not vary from the SAT implementation to ensure a fair comparison.

All models were trained for 300 epochs without early stopping. CIFAR-10, CIFAR-100, and SVHN were trained with a mini-batch size of 128 for training and 200 for testing. Due to resource limitations, Food-101 was trained with a mini-batch size of 128 for both training and testing.

The models were trained using SGD with an initial learning rate of 0.1 and a momentum of 0.9. The learning rate was reduced by 0.5 every 25 epochs. Weight decay was set to 0.0005.

For SAT and Socrates, an additional class (the idk class) was added, and the momentum of the loss was set to 0.9.

For Focal and Socrates, the gamma of the losses was set to 2, and alpha was set to 1.

The Selective Classification problem was evaluated with the coverage levels: 100, 98, 97, 95, 90, 85, 80, 75, 70, 60, 50, 40, 30, 20, and 10.

## D.3 DATASETS

As indicated by Feng et al. (2023), SAT was tested on easy-to-classify datasets. Therefore, for a more comprehensive analysis, we have selected a wide range of datasets with variable degrees of complexity. First, we chose the easy-to-classify CIFAR-10 and SVHN datasets. Although improvements may be less apparent with these toy datasets, the drawbacks of the methods could become more noticeable. To increase the challenge, we included the CIFAR-100 and Food-101 datasets. In particular, Food-101 serves as a good example of a real-world dataset, testing the reliability aspect of our new method. To further explore reliability, we tested the robustness of CIFAR-10 and CIFAR-100 using the CIFAR-10C and CIFAR-100C datasets as test sets.

The Street View House Number (SVHN) (Netzer et al., 2011) contains 73257 training and 26032 evaluation real-world small images of 32x32x3 with 10 classes. CIFAR-10 (Krizhevsky, 2009) comprises 50000 training and 10000 evaluation small images of 32x32x3 with 10 classes. CIFAR-100 (Krizhevsky, 2009) is like CIFAR-10 with 50000 training and 10000 evaluation small images of 32x32x3 but with 100 classes. CIFAR-10C (Hendrycks & Dietterich, 2019) comprises 50000 test small images of 32x32x3 with 10 classes created using the 10000 evaluation images using 5 different levels of corruption. CIFAR-100C (Hendrycks & Dietterich, 2019) similar to CIFAR-10C but with 100 classes. Food-101 (Bossard et al., 2014) constitutes 75750 training and 25250 evaluation images of 224x224x3 with 101 food classes.

## E SELECTIVE CLASSIFICATION ERROR RATE RESULTS AND ECE IN THE 300 EPOCH

As indicated in Subsection 6.2.1 of the main paper, the goal of this study is to produce calibrated selective classifiers that aim to achieve Selective Classification results similar to or better than those of SAT and CCL-SC, while ensuring well-calibrated confidence levels. The Selective Classification error rates achieved are comparable to or superior to those reached by SAT and CCL-SC. Notably, for the challenging CIFAR-100 and Food-101 datasets, Socrates significantly outperforms the Selective Classification error rates achieved by SAT. In this framework, once the model has been trained, it is insufficient to evaluate only the Selective Classification error rate without also considering metrics such as the ECE and accuracy.

The mean and standard deviation of the ECE values for the 300 epoch can be seen in Table 4 for the VGG-16 architecture and in Table 5 for the ResNet-34 architecture. The Selective Classification Error rates can be seen in Table 6 for the VGG-16 architecture, and in Table 7 for the Resnet-34 architecture.

Table 4: **ECE** values in a range of  $[0, 1]$  and **accuracy** (acc) values (100%) on the 300 epoch with the CIFAR-10, CIFAR-100, SVHN, Food-101, CIFAR-10C, and CIFAR-100 datasets with mean and standard deviation for trainings with **VGG-16 architecture**. A notable improvement can be seen in Food-101 dataset. Underline indicate the overall best performance, while bold highlight the best performance in each case.

Dataset	Coverage	<i>end-to-end case</i>		<i>first-epochs case</i>		
		Socrates (ours)	SAT	Socrates + Focal	CCL-SC + c.e	SAT + c.e
CIFAR-10	Acc Train	<b>97.47 ± 0.11</b>	94.33 ± 3.65	97.50 ± 0.28	<u>97.79 ± 0.11</u>	95.60 ± 3.80
	ECE Train	<b>0.003 ± 0.0004</b>	0.03 ± 0.01	<b>0.004 ± 0.0005</b>	0.007 ± 0.001	0.04 ± 0.008
	Acc val	<b>99.53 ± 0.03</b>	97.37 ± 2.36	99.81 ± 0.06	<u>99.93 ± 0.02</u>	98.50 ± 2.73
	ECE Val	<b>0.003 ± 0.0003</b>	0.02 ± 0.01	<b>0.004 ± 0.001</b>	0.005 ± 0.001	0.02 ± 0.001
	ECE Test	0.04 ± 0.002	<b>0.02 ± 0.005</b>	0.04 ± 0.002	<u>0.04 ± 0.001</u>	<b>0.02 ± 0.01</b>
CIFAR-100	Acc Train	<b>83.64 ± 0.30</b>	50.84 ± 1.09	84.18 ± 0.56	<u>85.59 ± 0.54</u>	69.93 ± 0.47
	ECE Train	<b>0.015 ± 0.001</b>	0.37 ± 0.015	0.03 ± 0.002	<u>0.017 ± 0.001</u>	0.10 ± 0.002
	Acc Val	<b>94.06 ± 0.18</b>	57.49 ± 1.37	95.74 ± 0.35	<u>97.04 ± 0.28</u>	88.38 ± 0.53
	ECE Val	<b>0.006 ± 0.001</b>	0.35 ± 0.02	<b>0.01 ± 0.001</b>	0.02 ± 0.001	0.05 ± 0.002
	ECE Test	<b>0.126 ± 0.004</b>	0.41 ± 0.01	<b>0.12 ± 0.003</b>	0.13 ± 0.002	0.14 ± 0.002
SVHN	Acc Train	<b>98.59 ± 0.1</b>	97.78 ± 0.04	<u>98.69 ± 0.1</u>	82.79 ± 35.67	78.81 ± 44.06
	ECE Train	<b>0.003 ± 0.0001</b>	0.01 ± 0.0001	<b>0.003 ± 0.0002</b>	0.004 ± 0.003	0.18 ± 0.36
	Acc Val	<b>99.42 ± 0.04</b>	98.82 ± 0.03	<u>99.64 ± 0.08</u>	85.57 ± 36.11	79.62 ± 44.51
	ECE Val	<b>0.002 ± 0.0002</b>	0.008 ± 0.0005	<b>0.002 ± 0.0001</b>	0.003 ± 0.002	0.17 ± 0.37
	ECE Test	0.013 ± 0.001	<b>0.007 ± 0.001</b>	0.015 ± 0.001	<b>0.012 ± 0.01</b>	0.17 ± 0.37
Food-101	Acc Train	<b>66.58 ± 0.86</b>	21.94 ± 1.78	66.51 ± 0.34	<u>68.98 ± 0.44</u>	40.06 ± 0.68
	ECE Train	<b>0.04 ± 0.002</b>	0.66 ± 0.03	<b>0.04 ± 0.003</b>	<u>0.04 ± 0.003</u>	0.31 ± 0.007
	Acc Val	<b>74.48 ± 0.74</b>	26.52 ± 2.15	74.60 ± 0.20	<u>75.58 ± 0.35</u>	55.08 ± 0.58
	ECE Val	<b>0.017 ± 0.002</b>	0.61 ± 0.03	<b>0.025 ± 0.004</b>	0.027 ± 0.004	0.20 ± 0.005
	ECE Test	<b>0.016 ± 0.003</b>	0.61 ± 0.03	0.017 ± 0.002	<b>0.011 ± 0.0003</b>	0.20 ± 0.01
CIFAR-10C	ECE Test	0.145 ± 0.002	<b>0.114 ± 0.003</b>	<u>0.154 ± 0.003</u>	0.156 ± 0.01	<b>0.11 ± 0.004</b>
CIFAR-100C	ECE Test	<b>0.24 ± 0.003</b>	0.51 ± 0.01	<b>0.24 ± 0.004</b>	0.25 ± 0.12	0.28 ± 0.001

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

Table 5: **ECE** values in a range of  $[0, 1]$  and **accuracy** (acc) values (100%) on the 300 epoch with the CIFAR-10, CIFAR-100, SVHN, Food-101, CIFAR-10C, and CIFAR-100 datasets with mean and standard deviation for trainings with **ResNet-34 architecture**. A notable improvement can be seen in Food-101 dataset. Underline indicate the overall best performance, while bold highlight the best performance in each case.

Dataset	Coverage	<i>end-to-end case</i>		<i>first-epochs case</i>		
		Socrates (ours)	SAT	Socrates + Focal	CCL-SC + c.e	SAT + c.e
CIFAR-10	Acc Train	<b>99.998 ± 0.001</b>	99.98 ± 0.02	<u>99.999 ± 0.001</u>	99.998 ± 0.002	99.994 ± 0.01
	ECE Train	<b>0.001 ± 0.0001</b>	0.002 ± 0.0002	<u>0.001 ± 0.00004</u>	0.001 ± 0.002	0.003 ± 0.0003
	Acc Val	<b>100 ± 0</b>	99.98 ± 0.003	<u>100 ± 0</u>	<b>100 ± 0</b>	99.99 ± 0.01
	ECE Val	<b>0.00077 ± 0.0001</b>	0.002 ± 0.0002	<u>0.00076 ± 0.0001</u>	0.0008 ± 0.0001	0.002 ± 0.0004
	ECE Test	<b>0.033 ± 0.0009</b>	0.033 ± 0.001	0.037 ± 0.003	<u>0.034 ± 0.002</u>	<b>0.032 ± 0.003</b>
CIFAR-100	Acc Train	<b>99.983 ± 0.006</b>	99.26 ± 0.07	99.97 ± 0.005	<b>99.9827 ± 0.007</b>	99.96 ± 0.01
	ECE Train	0.0061 ± 0.0006	<b>0.0058 ± 0.001</b>	0.0064 ± 0.0002	<b>0.006 ± 0.0004</b>	0.01 ± 0.001
	Acc Val	<b>99.9813 ± 0.005</b>	99.27 ± 0.07	<u>99.984 ± 0.003</u>	99.9836 ± 0.003	99.97 ± 0.01
	ECE Val	0.002 ± 0.0002	<b>0.001 ± 0.0003</b>	0.0025 ± 0.0002	<u>0.0002 ± 0.0001</u>	0.006 ± 0.0004
	ECE Test	<b>0.067 ± 0.01</b>	0.07 ± 0.02	0.07 ± 0.01	<u>0.064 ± 0.01</u>	0.065 ± 0.01
SVHN	Acc Train	<b>99.99 ± 0.002</b>	99.86 ± 0.02	<u>99.9948 ± 0.002</u>	99.99 ± 0.004	99.99 ± 0.002
	ECE Train	<b>0.00102 ± 0.0001</b>	0.002 ± 0.0004	<u>0.00104 ± 0.0001</u>	0.003 ± 0.001	0.002 ± 0.0002
	Acc Val	<b>99.997 ± 0.001</b>	99.86 ± 0.02	99.9957 ± 0.001	<b>99.996 ± 0.001</b>	99.995 ± 0.002
	ECE Val	<b>0.0008 ± 0.0001</b>	0.001 ± 0.0003	0.00067 ± 0.0001	<u>0.00065 ± 0.0001</u>	0.002 ± 0.0003
	ECE Test	<b>0.019 ± 0.001</b>	<b>0.019 ± 0.001</b>	0.021 ± 0.001	<u>0.02 ± 0.001</u>	<b>0.018 ± 0.001</b>
Food-101	Acc Train	<b>95.78 ± 0.31</b>	0 ± 0	82.38 ± 29.08	<b>95.52 ± 0.27</b>	88.85 ± 2.33
	ECE Train	<b>0.021 ± 0.002</b>	1 ± 0	0.026 ± 0.002	<b>0.023 ± 0.002</b>	0.08 ± 0.004
	Acc Val	<b>98.05 ± 0.37</b>	0 ± 0	82.95 ± 33.27	<b>97.57 ± 0.69</b>	92.19 ± 2.51
	ECE Val	<b>0.023 ± 0.004</b>	1 ± 0	0.04 ± 0.03	<b>0.026 ± 0.003</b>	0.06 ± 0.005
	ECE Test	<b>0.067 ± 0.001</b>	1 ± 0	<b>0.07 ± 0.02</b>	0.078 ± 0.002	0.09 ± 0.01
CIFAR10C	ECE Test	0.19 ± 0.01	<b>0.18 ± 0.01</b>	<u>0.19 ± 0.01</u>	0.18 ± 0.01	<b>0.17 ± 0.01</b>
CIFAR100C	ECE Test	<b>0.16 ± 0.04</b>	0.18 ± 0.04	<b>0.17 ± 0.03</b>	0.18 ± 0.02	0.20 ± 0.01

Table 6: **Selective Classification error rate** % on the 300 epoch with the CIFAR-10, CIFAR-100, SVHN, Food-101, CIFAR-10C, and CIFAR-100 datasets for various coverage rates % with mean and standard deviation for trainings with **VGG-16 architecture**. A notable improvement can be seen in Food-101 dataset. CCL-SC was not able to perform correctly for SVHN, SAT was not able for Food-101. Underline indicate the overall best performance, while bold highlight the best performance in each case.

Dataset	Coverage	<i>end-to-end case</i>		<i>first-epochs case</i>		
		Socrates (ours)	SAT	Socrates + Focal	CCL-SC + c.e	SAT + c.e
CIFAR-10	100	<b>6.44 ± 0.18</b>	7.08 ± 1.07	6.67 ± 0.19	<u>6.38 ± 0.14</u>	6.87 ± 1.08
	95	<b>4.14 ± 0.12</b>	4.78 ± 0.98	4.45 ± 0.14	<u>4.02 ± 0.14</u>	4.58 ± 1.12
	90	<b>2.43 ± 0.09</b>	3.01 ± 0.88	2.76 ± 0.13	<u>2.36 ± 0.13</u>	2.92 ± 1.01
	85	<b>1.48 ± 0.11</b>	1.82 ± 0.65	1.64 ± 0.20	<u>1.47 ± 0.16</u>	1.75 ± 0.74
	80	<b>0.85 ± 0.03</b>	1.12 ± 0.51	<b>1.05 ± 0.11</b>	1.18 ± 0.25	1.05 ± 0.46
	75	<u>0.52 ± 0.03</u>	0.67 ± 0.32	<i>0.68 ± 0.07</i>	1.05 ± 0.19	<b>0.61 ± 0.27</b>
	70	<b>0.38 ± 0.04</b>	0.43 ± 0.24	<i>0.51 ± 0.05</i>	0.91 ± 0.11	<b>0.42 ± 0.20</b>
CIFAR-100	100	<b>28.04 ± 0.24</b>	47.74 ± 1.27	28.08 ± 0.24	<u>28.01 ± 0.27</u>	<b>28.00 ± 0.06</b>
	95	<b>25.45 ± 0.32</b>	45.03 ± 1.33	<i>25.45 ± 0.32</i>	25.49 ± 0.31	<b>25.16 ± 0.08</b>
	90	<b>22.85 ± 0.30</b>	42.07 ± 1.38	23.07 ± 0.36	22.95 ± 0.28	<b>22.57 ± 0.07</b>
	85	<b>20.23 ± 0.30</b>	38.89 ± 1.45	20.76 ± 0.37	<i>20.37 ± 0.30</i>	<b>20.06 ± 0.08</b>
	80	<b>17.70 ± 0.23</b>	35.50 ± 1.43	18.34 ± 0.37	<i>17.79 ± 0.34</i>	<b>17.65 ± 0.09</b>
	75	<b>15.25 ± 0.28</b>	31.81 ± 1.48	15.85 ± 0.37	<i>15.27 ± 0.35</i>	<b>15.20 ± 0.12</b>
SVHN	100	2.72 ± 0.07	<b>2.65 ± 0.04</b>	<b>2.80 ± 0.03</b>	18.29 ± 34.73	18.22 ± 34.77
	95	1.15 ± 0.04	<u>1.04 ± 0.02</u>	<b>1.20 ± 0.08</b>	16.99 ± 35.46	16.89 ± 35.51
	90	0.74 ± 0.05	<b>0.61 ± 0.05</b>	<b>0.80 ± 0.05</b>	16.76 ± 35.58	16.57 ± 35.69
	85	0.62 ± 0.02	<b>0.45 ± 0.04</b>	<b>0.62 ± 0.05</b>	16.70 ± 35.62	16.44 ± 35.76
	80	0.55 ± 0.03	<b>0.38 ± 0.02</b>	<b>0.54 ± 0.05</b>	16.66 ± 35.64	16.39 ± 35.79
	75	0.49 ± 0.05	<b>0.33 ± 0.02</b>	<b>0.51 ± 0.03</b>	16.64 ± 35.65	16.35 ± 35.81
Food-101	100	<b>26.93 ± 0.52</b>	68.23 ± 2.19	<b>27.08 ± 0.25</b>	27.18 ± 0.19	29.00 ± 0.27
	95	<b>24.62 ± 0.54</b>	66.56 ± 2.31	24.78 ± 0.22	<b>24.62 ± 0.19</b>	26.63 ± 0.23
	90	<b>22.19 ± 0.50</b>	64.74 ± 2.43	22.40 ± 0.27	<b>22.04 ± 0.18</b>	24.29 ± 0.27
	85	<b>19.75 ± 0.52</b>	62.72 ± 2.56	20.00 ± 0.23	<b>19.38 ± 0.20</b>	21.89 ± 0.25
	80	<b>17.16 ± 0.60</b>	60.48 ± 2.70	17.59 ± 0.22	<b>16.75 ± 0.17</b>	19.43 ± 0.25
	75	<b>14.64 ± 0.63</b>	57.97 ± 2.86	15.18 ± 0.19	<b>14.19 ± 0.20</b>	17.06 ± 0.21
	70	<b>12.16 ± 0.56</b>	55.13 ± 3.04	12.87 ± 0.24	<b>11.67 ± 0.24</b>	14.66 ± 0.29
CIFAR10C	100	21.91 ± 0.24	<b>21.67 ± 0.22</b>	22.43 ± 0.35	<b>22.05 ± 0.57</b>	22.53 ± 1.95
	95	19.29 ± 0.26	<b>19.10 ± 0.25</b>	19.95 ± 0.42	<b>19.67 ± 0.59</b>	20.01 ± 2.04
	90	16.89 ± 0.28	<b>16.71 ± 0.25</b>	17.65 ± 0.45	<b>17.31 ± 0.59</b>	17.65 ± 2.11
	85	14.61 ± 0.28	<b>14.44 ± 0.24</b>	15.42 ± 0.46	<b>14.98 ± 0.59</b>	15.40 ± 2.16
	80	12.40 ± 0.27	<b>12.30 ± 0.23</b>	13.23 ± 0.43	<b>12.70 ± 0.58</b>	13.26 ± 2.16
	75	10.29 ± 0.27	<b>10.27 ± 0.23</b>	11.12 ± 0.40	<b>10.55 ± 0.55</b>	11.23 ± 2.16
CIFAR100C	100	<b>49.57 ± 0.14</b>	60.79 ± 0.78	49.68 ± 0.09	<i>49.43 ± 0.31</i>	<b>49.03 ± 0.20</b>
	95	<b>47.66 ± 0.13</b>	58.83 ± 0.81	47.59 ± 0.14	<i>47.56 ± 0.32</i>	<b>46.90 ± 0.21</b>
	90	<b>45.63 ± 0.13</b>	56.77 ± 0.84	<i>45.47 ± 0.18</i>	45.57 ± 0.34	<b>44.70 ± 0.21</b>
	85	<b>43.47 ± 0.12</b>	54.58 ± 0.85	<i>43.31 ± 0.23</i>	43.47 ± 0.34	<b>42.40 ± 0.22</b>
	80	<b>41.15 ± 0.12</b>	52.23 ± 0.87	<i>41.09 ± 0.27</i>	41.24 ± 0.36	<b>40.00 ± 0.21</b>
	75	<b>38.66 ± 0.13</b>	49.70 ± 0.89	<i>38.73 ± 0.30</i>	38.85 ± 0.36	<b>37.49 ± 0.21</b>
70	<b>36.01 ± 0.12</b>	46.94 ± 0.91	<i>36.22 ± 0.34</i>	36.30 ± 0.36	<b>34.86 ± 0.22</b>	

Table 7: **Selective Classification error rate %** on the 300 epoch with the CIFAR-10, CIFAR-100, SVHN, Food-101, CIFAR-10C, and CIFAR-100 datasets for various coverage rates % with mean and standard deviation for trainings with **ResNet-34 architecture**. A notable improvement can be seen in Food-101 dataset. SAT was not able to perform correctly for Food-101. Underline indicate the overall best performance, while bold highlight the best performance in each case.

Dataset	Coverage	<i>end-to-end case</i>		<i>first-epochs case</i>		
		Socrates (ours)	SAT	Socrates + Focal	CCL-SC + c.e	SAT + c.e
CIFAR-10	100	<b><u>4.95 ± 0.19</u></b>	5.10 ± 0.32	5.15 ± 0.26	<i>5.07 ± 0.10</i>	<b>4.97 ± 0.14</b>
	95	<b><u>2.71 ± 0.19</u></b>	2.85 ± 0.29	2.95 ± 0.20	<i>2.87 ± 0.11</i>	<b>2.84 ± 0.14</b>
	90	<b><u>1.46 ± 0.13</u></b>	1.55 ± 0.24	1.65 ± 0.17	<b>1.53 ± 0.12</b>	1.57 ± 0.15
	85	<b><u>0.81 ± 0.11</u></b>	0.88 ± 0.09	1.08 ± 0.12	<i>0.90 ± 0.11</i>	<b>0.90 ± 0.08</b>
	80	<b><u>0.56 ± 0.09</u></b>	0.60 ± 0.09	0.88 ± 0.09	<i>0.66 ± 0.06</i>	<b>0.60 ± 0.09</b>
	75	0.46 ± 0.07	<b><u>0.43 ± 0.11</u></b>	0.79 ± 0.10	<i>0.47 ± 0.03</i>	<b>0.44 ± 0.10</b>
	70	0.40 ± 0.09	<b><u>0.30 ± 0.08</u></b>	0.73 ± 0.08	<i>0.39 ± 0.04</i>	<b>0.36 ± 0.07</b>
CIFAR-100	100	<b><u>22.74 ± 0.34</u></b>	23.26 ± 0.50	23.19 ± 0.51	<i>23.23 ± 0.67</i>	<b>22.85 ± 0.27</b>
	95	<b><u>20.24 ± 0.45</u></b>	20.51 ± 0.51	20.48 ± 0.43	<i>20.44 ± 0.68</i>	<b>20.06 ± 0.37</b>
	90	<b><u>17.62 ± 0.59</u></b>	17.81 ± 0.43	18.05 ± 0.35	<i>17.80 ± 0.63</i>	<b>17.53 ± 0.43</b>
	85	<b><u>15.15 ± 0.54</u></b>	15.30 ± 0.46	15.71 ± 0.28	<i>15.23 ± 0.64</i>	<b>15.23 ± 0.36</b>
	80	12.85 ± 0.60	<b><u>12.83 ± 0.44</u></b>	13.62 ± 0.16	<b>12.90 ± 0.70</b>	13.03 ± 0.26
	75	10.74 ± 0.56	<b><u>10.73 ± 0.39</u></b>	11.68 ± 0.24	<b>10.69 ± 0.63</b>	11.09 ± 0.24
	70	8.73 ± 0.59	<b><u>8.65 ± 0.40</u></b>	9.84 ± 0.39	<b>8.67 ± 0.57</b>	9.19 ± 0.24
SVHN	100	<b><u>2.66 ± 0.09</u></b>	2.77 ± 0.08	2.78 ± 0.08	<i>2.74 ± 0.06</i>	<b>2.73 ± 0.11</b>
	95	1.02 ± 0.03	<b><u>0.99 ± 0.06</u></b>	1.11 ± 0.03	<b>1.03 ± 0.05</b>	1.04 ± 0.05
	90	0.65 ± 0.06	<b><u>0.60 ± 0.03</u></b>	0.76 ± 0.05	<i>0.67 ± 0.04</i>	<b>0.65 ± 0.07</b>
	85	0.55 ± 0.06	<b><u>0.48 ± 0.03</u></b>	0.67 ± 0.07	<i>0.59 ± 0.07</i>	<b>0.54 ± 0.06</b>
	80	0.52 ± 0.04	<b><u>0.43 ± 0.03</u></b>	0.60 ± 0.09	<i>0.56 ± 0.08</i>	<b>0.48 ± 0.05</b>
	75	0.48 ± 0.05	<b><u>0.40 ± 0.03</u></b>	0.56 ± 0.08	<i>0.54 ± 0.08</i>	<b>0.44 ± 0.05</b>
	70	0.46 ± 0.05	<b><u>0.39 ± 0.03</u></b>	0.53 ± 0.07	<i>0.51 ± 0.07</i>	<b>0.43 ± 0.05</b>
Food-101	100	<b><u>21.40 ± 0.79</u></b>	100 ± 0.0	32.33 ± 22.32	<i>22.77 ± 0.90</i>	<b>22.08 ± 0.75</b>
	95	<b><u>18.95 ± 0.80</u></b>	100 ± 0.0	30.20 ± 23.10	<i>20.09 ± 0.92</i>	<b>20.02 ± 0.74</b>
	90	<b><u>16.54 ± 0.75</u></b>	100 ± 0.0	28.23 ± 23.92	<b>17.39 ± 0.91</b>	17.97 ± 0.74
	85	<b><u>14.32 ± 0.74</u></b>	100 ± 0.0	26.37 ± 23.92	<b>14.75 ± 0.92</b>	15.99 ± 0.72
	80	<b><u>12.30 ± 0.78</u></b>	100 ± 0.0	24.60 ± 25.11	<b>12.30 ± 0.94</b>	14.08 ± 0.67
	75	<b><u>10.32 ± 0.68</u></b>	100 ± 0.0	22.94 ± 25.57	<b>10.00 ± 0.81</b>	12.20 ± 0.64
	70	<b><u>8.54 ± 0.62</u></b>	100 ± 0.0	21.49 ± 25.97	<b>7.85 ± 0.70</b>	10.37 ± 0.60
CIFAR10C	100	24.64 ± 0.49	<b><u>24.30 ± 0.89</u></b>	24.50 ± 0.66	<b>24.28 ± 0.46</b>	24.73 ± 0.67
	95	22.08 ± 0.49	<b><u>21.74 ± 0.92</u></b>	21.97 ± 0.67	<b>21.93 ± 0.47</b>	22.17 ± 0.72
	90	19.65 ± 0.46	<b><u>19.30 ± 0.92</u></b>	19.57 ± 0.70	<b>19.57 ± 0.47</b>	19.73 ± 0.77
	85	17.27 ± 0.44	<b><u>16.95 ± 0.89</u></b>	17.24 ± 0.73	<b>17.18 ± 0.47</b>	17.36 ± 0.85
	80	14.95 ± 0.41	<b><u>14.66 ± 0.85</u></b>	14.95 ± 0.77	<b>14.83 ± 0.46</b>	15.05 ± 0.95
	75	12.68 ± 0.37	<b><u>12.46 ± 0.78</u></b>	12.74 ± 0.80	<b>12.53 ± 0.44</b>	12.80 ± 1.02
	70	<b><u>10.56 ± 0.31</u></b>	10.69 ± 1.08	10.68 ± 0.81	<i>10.35 ± 0.39</i>	<b>10.22 ± 0.69</b>
CIFAR100C	100	<b><u>49.14 ± 0.32</u></b>	49.29 ± 0.72	49.73 ± 0.59	<i>49.04 ± 0.59</i>	<b>48.83 ± 0.24</b>
	95	47.46 ± 0.30	<b><u>47.38 ± 0.72</u></b>	47.73 ± 0.63	<i>47.11 ± 0.62</i>	<b>46.83 ± 0.28</b>
	90	45.56 ± 0.29	<b><u>45.38 ± 0.74</u></b>	45.69 ± 0.67	<i>45.06 ± 0.64</i>	<b>44.75 ± 0.31</b>
	85	43.49 ± 0.29	<b><u>43.24 ± 0.74</u></b>	43.58 ± 0.71	<i>42.92 ± 0.66</i>	<b>42.56 ± 0.33</b>
	80	41.26 ± 0.28	<b><u>40.95 ± 0.77</u></b>	41.37 ± 0.75	<i>40.66 ± 0.67</i>	<b>40.26 ± 0.33</b>
	75	38.87 ± 0.28	<b><u>38.48 ± 0.80</u></b>	39.05 ± 0.80	<i>38.28 ± 0.67</i>	<b>37.86 ± 0.33</b>
	70	36.33 ± 0.26	<b><u>35.83 ± 0.81</u></b>	36.64 ± 0.85	<i>35.76 ± 0.66</i>	<b>35.35 ± 0.34</b>

## F IS SELF-ADAPTIVE TRAINING LOSS A CALIBRATION LOSS? DETAILED ANALYSIS

**Overfitting:** As stated by Mukhoti et al. (2020) and Guo et al. (2017), overfitting appears to be linked to miscalibration. Therefore, the first step towards evaluating SAT as a calibrator is to examine the accuracy and loss curves to reconfirm the alleviation of the overfitting issue (one of the claims of the SAT and Focal methods). These curves are presented in Figure 6 and 7. Focal loss consistently maintains the same trend and does not induce overfitting with any dataset or architecture, except for the Food-101 dataset, where minor overfitting occurs in the initial epochs. In contrast, SAT loss behaves differently and does not consistently prevent overfitting. SAT shows overfitting with the SVHN and Food-101 datasets with the VGG-16 architecture during the *first-epochs* case, and for the *end-to-end* case with the challenging CIFAR-100 and Food-101 datasets across both architectures. Furthermore, with the CIFAR-100 and Food-101 datasets with the VGG-16 architecture, the accuracy achieved with SAT loss is significantly lower than with Focal loss. For the *end-to-end* case with Food-101, and ResNet-34, SAT could not train. Although these observations suggest that SAT may not be an effective calibration loss, they do not provide definitive evidence of miscalibration, necessitating further analysis.

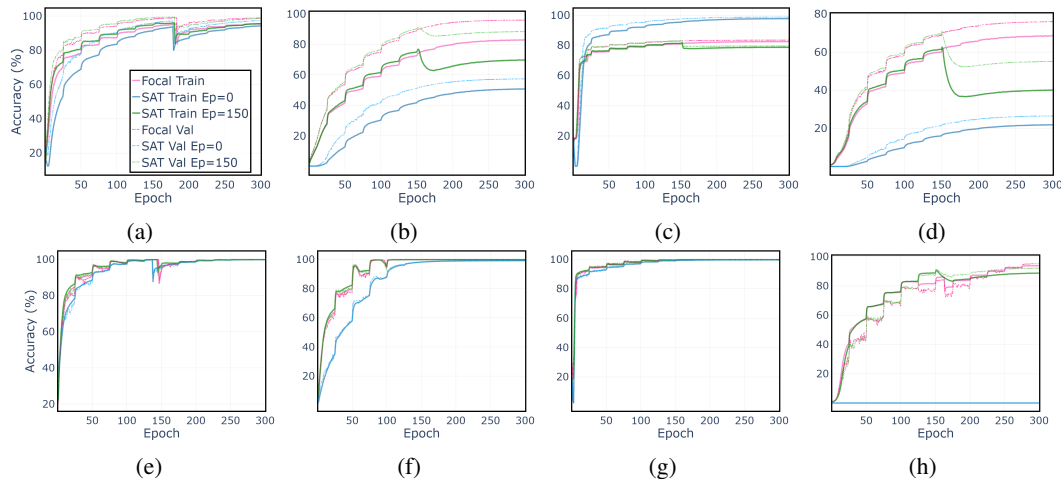


Figure 6: Accuracy curves of models trained on CIFAR-10 (a and e), CIFAR-100 (b and f) SVHN (c and g) and Food-101 (d and h) datasets using Focal and SAT (*first-epochs* and *end-to-end* cases) methods with VGG-16 (a, b, c, and d) and ResNet-34 (e, f, g, and h) architectures.

**Calibration Metrics:** The second step towards the analysis of SAT as a calibration loss is to visualize specific calibration metrics. The snapshot of the ECE and MCE metrics in the reliability diagram of the last training epoch does not give enough insights to output calibration conclusions, instead, guided for the experimentation phase made in Mukhoti et al. (2020), the ECE and MCE values in each epoch of the training process produce noticeable insights. Therefore, the visualization starts analyzing the ECE value along the epochs. These ECE values along epochs curves can be seen in Figure 8.

The ECE along epochs curves of Focal loss exhibit a consistent downward trend across both architectures, except for VGG-16 with CIFAR-100 and Food-101 datasets where there is an increase in the initial epochs but in an acceptable ECE range. Regarding SAT, in the *first-epochs* case, ECE values for VGG-16 architectures rise significantly after 150 epochs, especially in the Food-101 dataset, but this increase is less noticeable for ResNet-34 architectures. In the *end-to-end* case, both architectures show high initial ECE values that gradually decrease, though VGG-16 has a particularly high ECE of around 0.9, compared to much lower values with focal loss. These observations suggest that SAT is less reliable as a calibrator compared to focal loss, which performs more consistently.

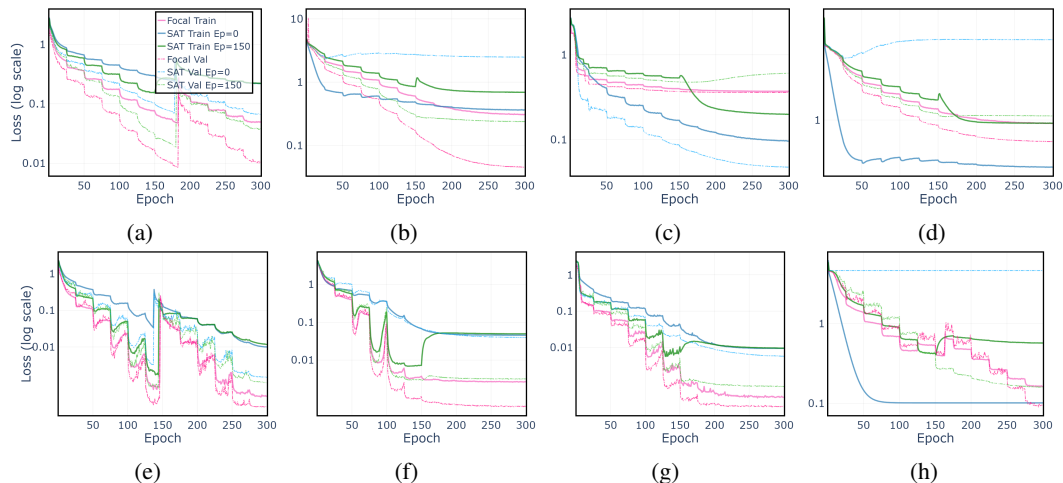


Figure 7: Loss curves of models trained on CIFAR-10 (a and e), CIFAR-100 (b and f), SVHN (c and g), and Food-101 (d and h) datasets using Focal and SAT (*first-epochs* and *end-to-end* cases) methods with VGG-16 (a, b, c, and d), and ResNet-34 (e, f, g, and h) architectures.

The MCE curves do not provide sufficient insights, as the MCE values are typically driven by only a few instances. The main claim is that the SAT exhibits distinct trends in both VGG-16 and ResNet-34 architectures compared to Focal loss, mirroring the observations made regarding ECE.

Therefore, a significant claim can be put forth: SAT loss does not seem to be a good loss for training calibrated models, and it appears detrimental when the goal is to train for a small number of epochs. It is well-known that the aim is not always to train for longer, as it depends on the dataset and the architecture. In this case, the SAT loss outputs calibrated confidences after a considerable amount of epochs, which may not be desirable in all cases.

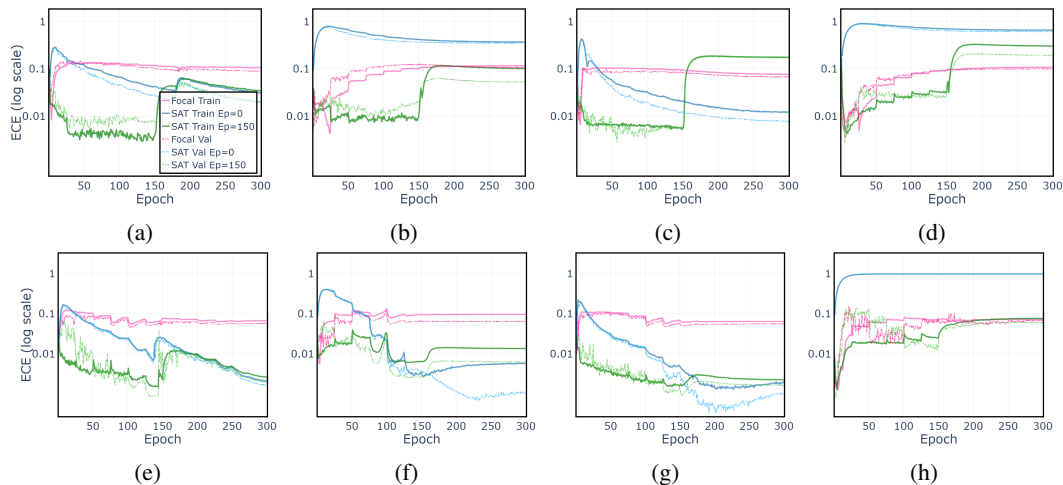


Figure 8: Evolution of the Expected Calibration Error (ECE) across epochs for models trained on CIFAR-10 (a and e), CIFAR-100 (b and f), SVHN (c and g), and Food-101 datasets (d and h) using Focal and Self-Adaptive Training (*first-epochs* and *end-to-end* cases) methods with VGG-16 (a, b, c and d), and ResNet-34 (e, f, g, h) architectures.

**Idk class:** Given that the additional idk class retains the model’s knowledge when it does not know, it is reasonable to anticipate that this class will change across epochs, typically exhibiting a decreasing trend. Analyzing the average values of the idk class confidences across epochs provides valuable insights; these plots are shown in Figure 9. Visualising these curves, the idk class appears to

be directly related to calibration. If we compare the ECE across epochs curves with the average of the idk confidences across epochs curves, it is noticeable that both values plot similar trends. When the model believes that it is more certain about what it does not know (indicated by higher average idk confidences), the ECE value tends to be larger, which could be possible due to incorrect confidence values associated with the ground truth classes. This assumption prompts us to consider: *Might the extra idk class approach be beneficial in some way in the calibration aspect or detrimental?* This visualized behaviour was the main source of inspiration to decide to add the predictions associated with the unknown class in the novel Socrates loss to calibrate the training.

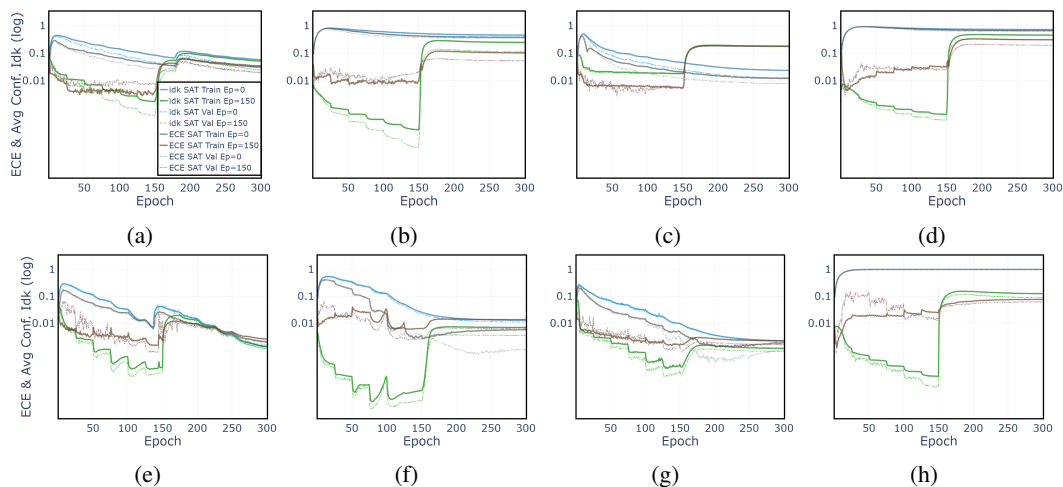


Figure 9: Curves depicting the average values of the idk class confidences across the epochs of models trained on CIFAR-10 (a and e), CIFAR-100 (b and f), SVHN (c and g), and Food-101 (d and h) datasets using Focal and SAT (*first-epochs* and *end-to-end* cases) methods with VGG-16 (a, b, c, and d), and ResNet-34 (e, f, g, and h) architectures.

**Self-Adaptive Training (SAT) loss seems not to be a calibration loss:** Based on the aforementioned empirical analysis the following claim can be made: Unlike Focal loss, which produces very well-calibrated models and follows similar trends across all the datasets and architectures, SAT loss exhibits certain tendencies that ultimately lead to the conclusion that it is not a loss that allows learning calibrated models in all the scenarios, especially when aiming to train for a small number of epochs or when dealing with complex datasets such as Food-101. Additionally, when the loss is used *end-to-end*, the miscalibration in the first epochs is excessively large, and in some cases (CIFAR-100 and Food-101 with VGG-16) it remains significantly large until the end of training. When the loss is applied after the initial epochs (*first-epochs* case), miscalibration begins to emerge.

## G SOCRATES LOSS AS A CALIBRATOR: FIGURES

Due to space constraints, the graphs for all datasets and architectures evaluating the calibration capacity of the Socrates method versus the CCL-SC method are presented in this section. This is supplementary material of Subsection 6.2.

**Accuracy and Loss curves:** The accuracy and loss curves have provided insightful visualizations of performance and fitting. The accuracy curves can be found in Figure 10 and the loss curves in Figure 11.

**ECE values across epochs curves:** The curves plotting the Expected Calibration Error (ECE) values across epochs serve as the focal point of this research, offering key insights into the calibration capacity of the methods. The ECE values across epochs curves can be found in Figure 12.

**Idk class:** Given that the additional idk class retains the model’s knowledge when it does not know, it is reasonable to anticipate that this class will change across epochs, typically exhibiting a



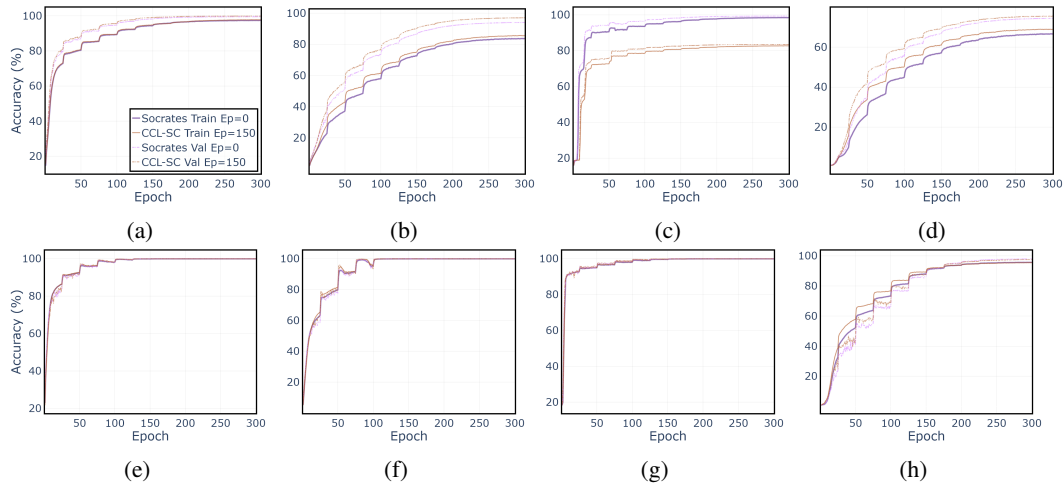


Figure 10: Accuracy across epochs curves of models trained on CIFAR-10 (a and e), CIFAR-100 (b and f), SVHN (c and g), and Food-101 (d and h) datasets using Socrates and CCL-SC methods with VGG-16 (a, b, c and d) and ResNet-34 (e, f, g, h) architectures.

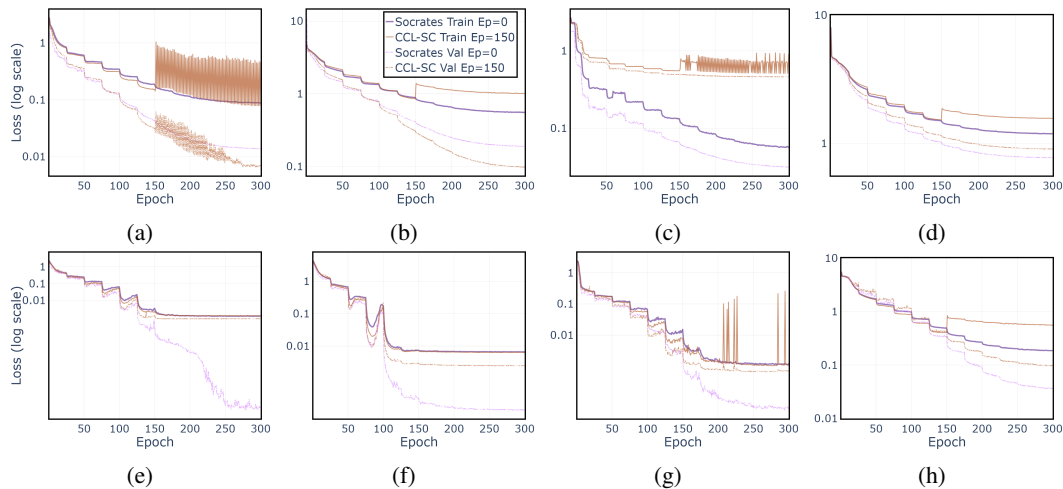


Figure 11: Loss curves of models trained on CIFAR-10 (a and e), CIFAR-100 (b and f), SVHN (c and g), and Food-101 (d and h) datasets using using Socrates and CCL-SC methods (*first-epochs* and *end-to-end* cases) with VGG-16 (a, b, c and d) and ResNet-34 (e, f, g, h) architectures.

decreasing trend. Analyzing the average values of the *idk* class confidences across epochs provides valuable insights; these plots are shown in Figure 13.

1350  
 1351  
 1352  
 1353  
 1354  
 1355  
 1356  
 1357  
 1358  
 1359  
 1360  
 1361  
 1362  
 1363  
 1364  
 1365  
 1366  
 1367  
 1368  
 1369  
 1370  
 1371  
 1372  
 1373  
 1374  
 1375  
 1376  
 1377  
 1378  
 1379  
 1380  
 1381  
 1382  
 1383  
 1384  
 1385  
 1386  
 1387  
 1388  
 1389  
 1390  
 1391  
 1392  
 1393  
 1394  
 1395  
 1396  
 1397  
 1398  
 1399  
 1400  
 1401  
 1402  
 1403

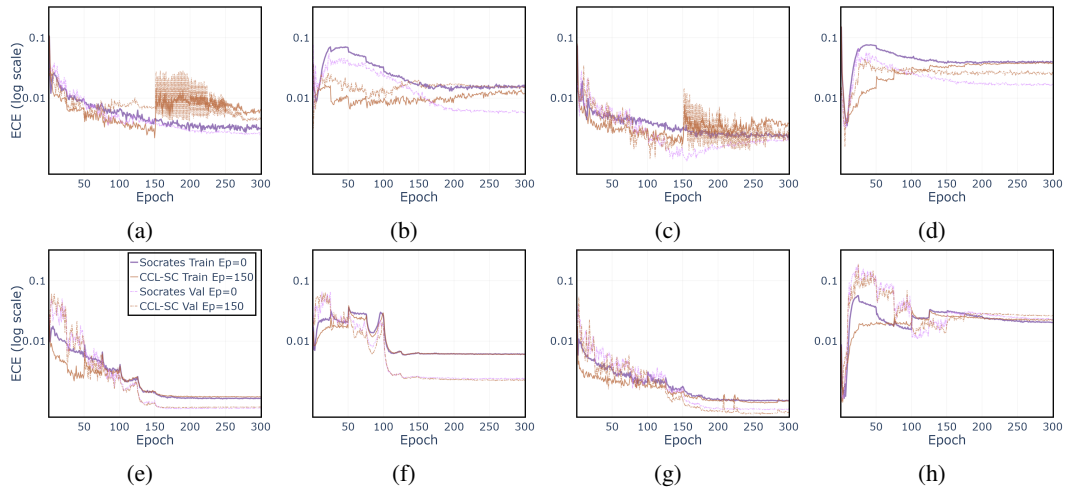


Figure 12: Evolution of the Expected Calibration Error (ECE) across epochs for models trained on CIFAR-10 (a and e), CIFAR-100 (b and f), SVHN (c and g) and Food-101 (d and h) datasets using Socrates and CCL-SC methods (*first-epochs* and *end-to-end* cases) with VGG-16 (a, b, c and d) and ResNet-34 (e, f, g, h) architectures.

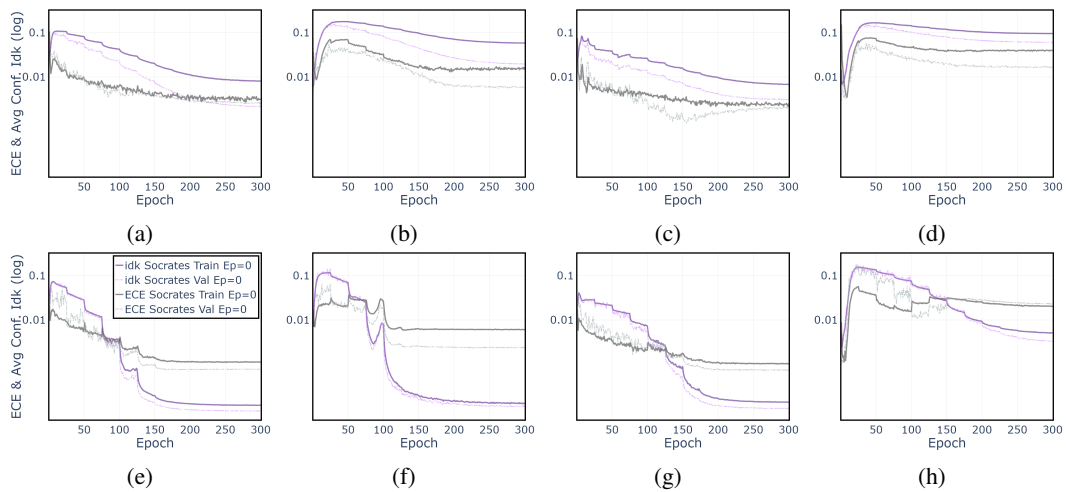


Figure 13: Curves depicting the average values of the idk class confidences across the epochs of models trained on CIFAR-10 (a and e), CIFAR-100 (b and f), SVHN (c and g), and Food-101 (d and h) datasets trained using Socrates and CCL-SC methods (*first-epochs* and *end-to-end* cases) with VGG-16 (a, b, c and d) and ResNet-34 (e, f, g, h) architectures.

## H RISK-COVERAGE CURVES

The risk-coverage curves offer a clear representation of the power of Socrates compared to the CCL-SC method. As shown in figure 14, these curves illustrate Socrates reaches similar values or outperforms, thereby providing a more reliable framework for model evaluation.

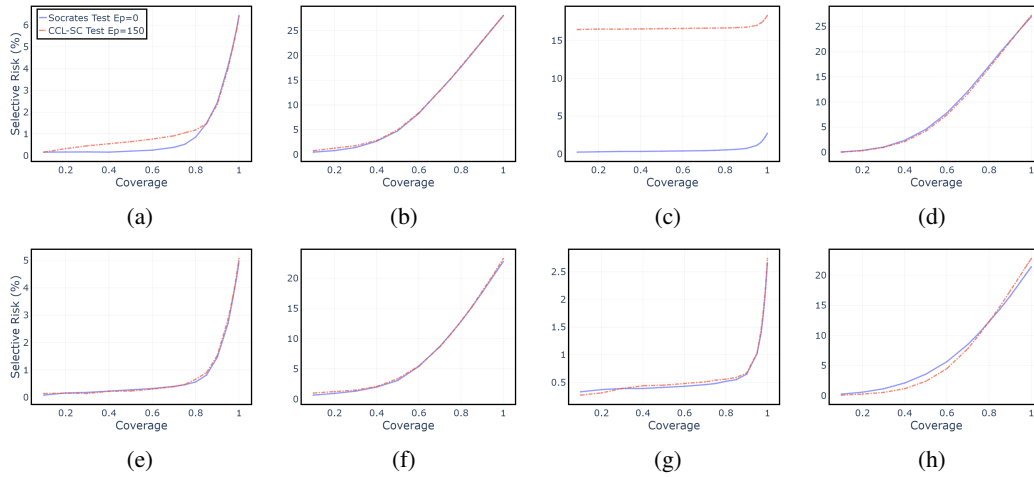


Figure 14: Risk-Coverage curves of models trained on CIFAR-10 (a and e), CIFAR-100 (b and f), SVHN (c and g), and Food-101 (d and h) datasets using Socrates (*end-to-end* case) and CCL-SC (*first-epochs* case) methods with VGG-16 (a,b, c and d) and ResNet-34 (e, f, g, h) architectures.