

## A RELATED WORKS

**Calibration Methods** Post-hoc calibration methods adjust model outputs after training to improve calibration. A widely used technique is Temperature Scaling (TS) (Guo et al., 2017), which smooths softmax probabilities by search a temperature factor on a validation set. Enhanced variants of TS include Parameterized Temperature Scaling (PTS) (Tomani et al., 2022), which uses a neural network to learn the temperature, and Class-based Temperature Scaling (CTS) (Frenkel et al., 2021), which applies adjustments on a class-wise basis. Group Calibration (GC) (Yang et al., 2024) and ProCal (Xiong et al., 2023a) aim for multi-calibration (Hébert-Johnson et al., 2018) by splitting data samples by proximity and grouping. Another stream of work is train-time calibration such as Brier Loss (Brier, 1950), Dirichlet Scaling (Kull et al., 2019), Maximum Mean Calibration Error (MMCE) (Kumar et al., 2018), Label Smoothing (Szegedy et al., 2016), and Focal Loss (Mukhoti et al., 2020) and Dual Focal Loss (Tao et al., 2023). However, these methods often require substantial higher computational overhead.

**Ensemble-Based Calibration** Ensemble-based methods ensemble multiple outputs in different ways. They use models or samples to approximate Bayesian Inference. Lakshminarayanan et al. (2017) propose deep ensembles as a scalable alternative to Bayesian Neural Networks (BNNs) for uncertainty estimation. Similarly, Gal & Ghahramani (2016) treat dropout as approximate Bayesian inference. Data-centric ensemble techniques using test-time augmentation, as described by Conde et al. (2023), also help improve calibration. Zhang et al. (2020) resort to the power of Bayesian inference and proposed a Ensemble-based TS (ETS). However, these methods typically require significant computational resources to train multiple models or perform repeated inferences. In contrast, our approach relies on consistency rather than probability distribution modeling.

**Consistency in LLMs** Consistency has emerged as a key approach for black-box uncertainty estimation and hallucination detection in large language models (LLMs). These methods evaluate uncertainty by measuring variability in outputs across slight changes, such as different sampling techniques or rephrased prompts. Confident models produce stable outputs, while variability indicates uncertainty. For instance, SelfCheckGPT (Manakul et al., 2023) uses sampling and similarity metrics like BERTScore and NLI to detect hallucinations, while Lin et al. (2023) analyze a similarity matrix to estimate uncertainty. Xiong et al. (2023b) further break down uncertainty estimation into prompting, sampling, and consistency-based aggregation. These methods, which rely on output stability, are efficient alternatives to probabilistic approaches.

## B PERTURBATION OF DIFFERENT LAYER

This section presents a detailed analysis of the impact of perturbations applied at various levels of a ResNet50 model, trained on CIFAR-10. The experiments were conducted using 32 samples, and the effects on ECE, accuracy, and optimal perturbation values were evaluated.

Perturbation Level	ECE (%)	Accuracy (%)	Optimal Perturbation
Image	1.1	95.25	train aug jitter0.1
Logits	0.73	95.04	8.2
Feature (Last Layer)	2.06	95.06	3.0
Feature (Layer 4)	0.53	95.29	13.28
Feature (Layer 3)	53.12	10.03	20.12
Feature (Layer 2)	56.28	10.02	20.21
Feature (Layer 1)	49.53	10.11	20.75

Table 5: Comparison of perturbations at different layers with number of samples set to 32 using ECE $\downarrow$  and Accuracy $\uparrow$ , evaluated on ResNet50 with CIFAR-10. ECE values are reported with 15 bins. Optimal Perturbations for logits and features are represented in  $\epsilon$  value

From Table 5, we observe a clear trend in the performance of perturbations applied at different layers of the model. Perturbation at the logits level achieves a favorable trade-off between calibration and efficiency. Although the perturbation applied to the fourth layer’s feature space slightly improves the ECE to 0.53%, the associated computational cost is significantly higher, with the optimal perturbation value of 13.28.

On the other hand, perturbations applied at lower feature levels (Layer 1 to Layer 3) result in severe degradation of both accuracy and calibration. Specifically, the ECE increases drastically to above 50%, and accuracy drops to approximately 10%, with a significant increase in computing time and memory use. This suggests that perturbing the features at these lower layers disrupts the model’s ability to recognize patterns and correctly classify the input data. We hypothesize that this is due to the higher sensitivity of lower layers to the raw data structure, where perturbations may significantly distort the features necessary for effective recognition.

## C COMPARISON OF POST-HOC CALIBRATION METHODS ON OTHER METRICS

As shown in table 6. The proposed CC method consistently achieves the lowest AdaECE values, outperforming the other methods. This indicates better calibration performance, in line with our discussion in the main text. For instance, in CIFAR-10, Wide-ResNet has an AdaECE of 0.40 with CC compared to 3.24 for Vanilla, showing a significant improvement. Similar results are observed across other models and datasets. The formula for Adaptive-ECE is as follows:

$$\text{Adaptive-ECE} = \sum_{i=1}^B \frac{|B_i|}{N} |I_i - C_i| \text{ s.t. } \forall i, j \cdot |B_i| = |B_j| \quad (11)$$

Dataset	Model	Vanilla	TS	ETS	PTS	CTS	GC	CC (ours)
CIFAR-10	ResNet-50	4.33	2.14	2.14	2.14	1.71	1.24	<b>0.64</b>
	ResNet-110	4.40	1.89	1.89	1.90	1.31	<b>0.94</b>	0.96
	DenseNet-121	4.49	2.12	2.12	2.12	1.71	1.28	<b>1.20</b>
	Wide-ResNet	3.24	1.71	1.71	1.71	1.42	1.17	<b>0.40</b>
CIFAR-100	ResNet-50	17.52	5.76	5.72	5.66	5.79	3.43	<b>1.61</b>
	Wide-ResNet	15.34	4.48	4.45	4.41	4.69	2.24	<b>1.73</b>
ImageNet	ResNet-50	3.73	2.07	2.07	2.06	3.22	2.56	<b>1.47</b>
	DenseNet-121	6.59	1.67	1.68	1.69	1.89	2.49	<b>1.36</b>
	Wide-ResNet-50	5.32	2.97	2.97	2.95	4.13	2.18	<b>1.27</b>
	ViT-B-16	5.59	4.05	4.06	4.08	5.50	1.86	<b>1.76</b>
	ViT-B-32	6.40	3.83	3.85	3.91	5.73	<b>1.33</b>	1.77

Table 6: **Comparison of Post-Hoc Calibration Methods Using AdaECE $\downarrow$  Across Various Datasets and Models.** AdaECE values are reported with 15 bins. The best results for each combination is in bold, and our method (CC) is highlighted. Results are averaged over 5 runs.

As shown in table 7. The CC method also performs the best in terms of class-wise calibration, with consistently lower CECE values. This confirms that CC provides better calibration across individual classes, as discussed in the main body. For example, for ResNet-50 on CIFAR-100, CC achieves a CECE of 0.20, which is the lowest among the methods. CECE is another measure of calibration performance that addresses the deficiency of ECE in only measuring the calibration performance of the single predicted class. It can be formulated as:

$$\text{Classwise-ECE} = \frac{1}{\mathcal{K}} \sum_{i=1}^B \sum_{j=1}^{\mathcal{K}} \frac{|B_{i,j}|}{N} |I_{i,j} - C_{i,j}| \quad (12)$$

As shown in table 8, interestingly, the NLL values are generally higher with the CC method compared to some other calibration methods, despite its superior calibration performance in AdaECE and CECE. This suggests that while CC improves calibration, it may come at the cost of slightly higher NLL values. For instance, for CIFAR-100 on ResNet-50, CC has a higher NLL than TS, but it remains competitive overall.

9 indicates that there is little to no change in accuracy across the calibration methods, with all methods performing similarly in terms of classification accuracy. This pattern is consistent with the main section, showing CC improves calibration without sacrificing accuracy. For example, on CIFAR-10, Wide-ResNet achieves almost identical accuracy for all methods, with CC slightly outperforming others in specific cases.

Dataset	Model	Vanilla	TS	ETS	PTS	CTS	GC	CC (ours)
CIFAR-10	ResNet-50	0.91	0.45	0.45	0.45	0.41	0.46	<b>0.39</b>
	ResNet-110	0.92	0.48	0.48	0.48	0.42	0.52	<b>0.41</b>
	DenseNet-121	0.92	0.48	0.48	0.48	<b>0.41</b>	0.54	0.43
	Wide-ResNet	0.68	0.37	0.37	0.37	0.37	0.48	<b>0.32</b>
CIFAR-100	ResNet-50	0.38	0.21	0.21	0.21	0.22	0.21	<b>0.20</b>
	Wide-ResNet	0.34	0.19	0.19	0.19	0.20	0.20	<b>0.18</b>
ImageNet	ResNet-50	<b>0.03</b>	<b>0.03</b>	<b>0.03</b>	<b>0.03</b>	<b>0.03</b>	<b>0.03</b>	<b>0.03</b>
	DenseNet-121	<b>0.03</b>	<b>0.03</b>	<b>0.03</b>	<b>0.03</b>	<b>0.03</b>	<b>0.03</b>	<b>0.03</b>
	Wide-ResNet-50	0.03	0.03	0.03	0.03	0.03	0.03	<b>0.02</b>
	ViT-B-16	0.03	<b>0.02</b>	<b>0.02</b>	<b>0.02</b>	0.03	0.02	<b>0.02</b>
	ViT-B-32	<b>0.03</b>	<b>0.03</b>	<b>0.03</b>	<b>0.03</b>	<b>0.03</b>	<b>0.03</b>	<b>0.03</b>

Table 7: Comparison of Post-Hoc Calibration Methods Using CECE↓ Across Various Datasets and Models. CECE values are reported with 15 bins. The best-performing method for each dataset-model combination is in bold, and our method (CC) is highlighted. Results are averaged over 5 runs.

Dataset	Model	Vanilla	TS	ETS	PTS	CTS	GC	CC (ours)
CIFAR-10	ResNet-50	41.21	20.39	20.39	20.38	20.15	<b>19.97</b>	20.39
	ResNet-110	47.52	21.52	21.52	21.52	20.84	<b>20.68</b>	23.33
	DenseNet-121	42.93	21.78	21.78	21.78	21.01	<b>20.30</b>	22.19
	Wide-ResNet	26.75	15.33	15.33	15.33	<b>15.13</b>	15.32	17.10
CIFAR-100	ResNet-50	153.67	<b>106.07</b>	<b>106.07</b>	<b>106.07</b>	106.25	107.80	108.40
	Wide-ResNet	140.11	<b>95.71</b>	<b>95.71</b>	<b>95.71</b>	96.38	96.92	99.30
ImageNet	ResNet-50	96.12	94.82	94.82	<b>94.81</b>	99.58	99.07	140.57
	DenseNet-121	109.52	<b>103.90</b>	<b>103.90</b>	103.91	106.13	108.14	162.02
	Wide-ResNet-50	88.56	<b>86.46</b>	<b>86.46</b>	<b>86.46</b>	91.68	nan	120.59
	ViT-B-16	83.71	78.63	78.63	<b>78.63</b>	85.19	82.14	106.89
	ViT-B-32	107.76	101.67	101.67	<b>101.66</b>	107.53	105.45	141.71

Table 8: Comparison of Post-Hoc Calibration Methods Using NLL↓ Across Various Datasets and Models. The best-performing method for each dataset-model combination is in bold, and our method (CC) is highlighted. Results are averaged over 5 runs.

In figure 6 we see that the proposed CC method significantly reduces both AdaECE and CECE values compared to other calibration methods, indicating better calibration for Wide-ResNet on CIFAR-10. The accuracy remains mostly unchanged across all methods, while NLL is slightly higher for CC compared to other methods like TS and ETS. This behavior is consistent with our findings in the main text.

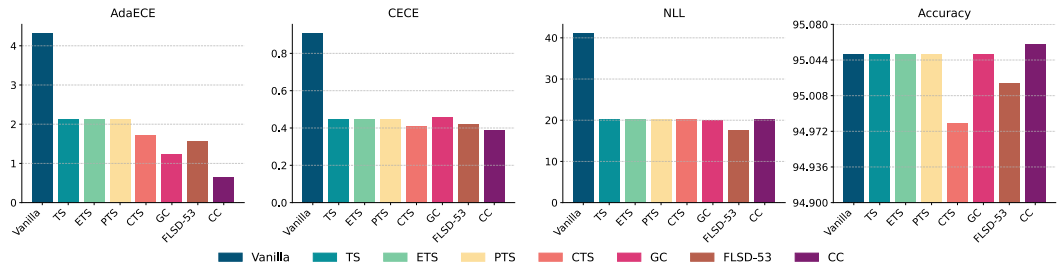


Figure 6: Calibration performance of ResNet-50 on Cifar-10 using AdaECE↓, CECE↓, NLL↓, and Accuracy↑. ECE, AdaECE, and CECE are reported with 15 bins. Colors in the legend represent different methods. Results are averaged over 5 runs.

In Figure 6 for ResNet-50 on CIFAR-10, the CC method demonstrates excellent performance with the lowest AdaECE and CECE values, further supporting its effectiveness in calibration. NLL is higher for CC, which is interesting given its superior performance in other metrics. However, accuracy remains largely unchanged, consistent with the overall findings discussed in the text.

Figure 8 illustrates the performance of ResNet-50 on CIFAR-100 across different calibration methods. The proposed CC method again shows the lowest AdaECE and CECE, confirming its superior

Dataset	Model	Vanilla	TS	ETS	PTS	CTS	GC	CC (ours)
CIFAR-10	ResNet-50	95.05	95.05	95.05	95.05	94.98	95.05	<b>95.06</b>
	ResNet-110	95.11	95.11	95.11	95.11	<b>95.18</b>	95.11	95.16
	DenseNet-121	95.02	95.02	95.02	95.02	95.01	95.02	<b>95.04</b>
	Wide-ResNet	<b>96.13</b>	<b>96.13</b>	<b>96.13</b>	<b>96.13</b>	96.06	<b>96.13</b>	<b>96.13</b>
CIFAR-100	ResNet-50	76.70	76.70	76.70	76.70	<b>76.72</b>	76.70	76.71
	Wide-ResNet	79.29	79.29	79.29	79.29	79.17	79.29	<b>79.31</b>
ImageNet	ResNet-50	<b>76.08</b>	<b>76.08</b>	<b>76.08</b>	<b>76.08</b>	74.62	<b>76.08</b>	<b>76.08</b>
	DenseNet-121	74.16	74.16	74.16	74.16	73.08	74.16	<b>74.37</b>
	Wide-ResNet-50	78.40	78.40	78.40	78.40	77.07	78.40	<b>78.48</b>
	ViT-B-16	<b>81.09</b>	<b>81.09</b>	<b>81.09</b>	<b>81.09</b>	80.01	<b>81.09</b>	81.06
	ViT-B-32	<b>75.94</b>	<b>75.94</b>	<b>75.94</b>	<b>75.94</b>	74.90	<b>75.94</b>	75.90

Table 9: Comparison of Post-Hoc Calibration Methods Using Accuracy $\uparrow$  Across Various Datasets and Models. Top-1 accuracy values are reported. The best results for each combination is in bold, and our method (CC) is highlighted. Results are averaged over 5 runs.

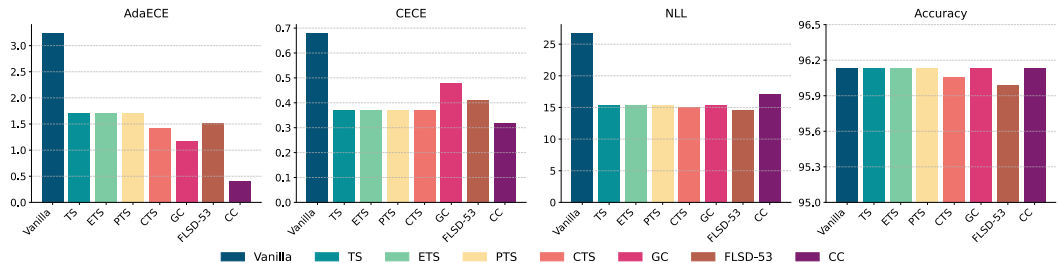


Figure 7: Calibration performance of Wide-ResNet on CIFAR-10 using AdaECE $\downarrow$ , CECE $\downarrow$ , NLL $\downarrow$ , and Accuracy $\uparrow$ . ECE, AdaECE, and CECE are reported with 15 bins. Colors in the legend represent different methods. Results are averaged over 5 runs.

calibration performance. NLL for CC is slightly higher compared to TS, but accuracy shows minimal changes across methods. These results align with our overall conclusions that CC improves calibration without sacrificing accuracy.

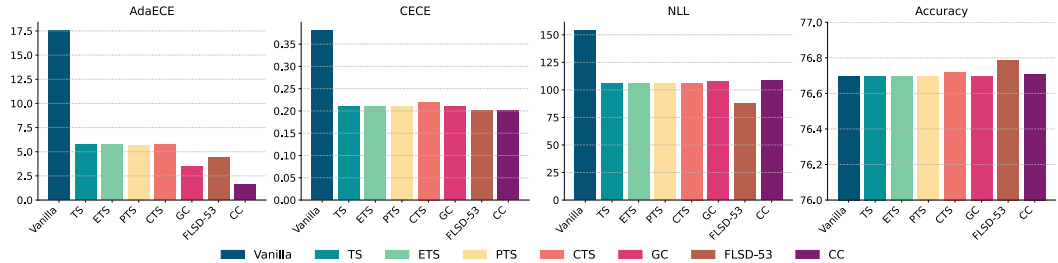


Figure 8: Calibration performance of ResNet-50 on CIFAR-100 using AdaECE $\downarrow$ , CECE $\downarrow$ , NLL $\downarrow$ , and Accuracy $\uparrow$ . ECE, AdaECE, and CECE are reported with 15 bins. Colors in the legend represent different methods. Results are averaged over 5 runs.

## D COMPARISON OF VARIOUS TRAINING-TIME CALIBRATION METHODS ON OTHER METRICS

As shown in Table 10, CC consistently outperforms baseline models across all metrics and datasets. Specifically, on CIFAR-10 and CIFAR-100, CC achieves significantly lower AdaECE scores for ResNet-50, ResNet-110, DenseNet-121, and Wide-ResNet compared to traditional methods such as Brier Loss, and MMCE. For instance, on CIFAR-100 with ResNet-110, CC reduces the AdaECE from 19.05 (baseline) to 5.28, showing superior calibration performance.

Dataset	Model	Cross-Entropy		Brier Loss		MMCE		LS-0.05		FLSD-53		FL-3	
		base	ours	base	ours	base	ours	base	ours	base	ours	base	ours
CIFAR-10	ResNet-50	4.33	<b>0.64</b>	1.75	<b>0.99</b>	4.55	<b>1.06</b>	3.88	<b>1.74</b>	1.56	<b>0.36</b>	1.95	<b>0.71</b>
	ResNet-110	4.40	<b>0.96</b>	2.60	<b>0.30</b>	5.07	<b>1.80</b>	4.48	<b>2.43</b>	2.08	<b>0.73</b>	1.64	<b>0.38</b>
	DenseNet-121	4.49	<b>1.20</b>	2.02	<b>0.64</b>	5.10	<b>1.76</b>	4.40	<b>1.94</b>	1.38	<b>0.53</b>	1.23	<b>0.69</b>
	Wide-ResNet	3.24	<b>0.40</b>	1.70	<b>0.57</b>	3.29	<b>0.63</b>	4.27	<b>1.54</b>	1.52	<b>0.42</b>	1.84	<b>0.42</b>
CIFAR-100	ResNet-50	17.52	<b>1.61</b>	6.55	<b>1.90</b>	15.32	<b>1.88</b>	7.66	<b>6.17</b>	4.39	<b>1.48</b>	5.09	<b>1.70</b>
	ResNet-110	19.05	<b>5.28</b>	7.72	<b>3.54</b>	19.14	<b>5.14</b>	11.14	<b>8.00</b>	8.56	<b>3.50</b>	8.64	<b>3.98</b>
	DenseNet-121	20.99	<b>5.85</b>	5.04	<b>2.02</b>	19.10	<b>3.90</b>	12.83	<b>7.06</b>	3.54	<b>1.52</b>	4.14	<b>2.03</b>
	Wide-ResNet	15.34	<b>1.73</b>	4.28	<b>1.92</b>	13.16	<b>2.06</b>	5.14	<b>4.75</b>	2.77	<b>1.79</b>	2.07	<b>1.58</b>

Table 10: Comparison of Train-time Calibration Methods Using AdaECE↓ Across Various Datasets and Models. AdaECE values are reported with 15 bins. The best results for each combination is in bold, and our method (CC) is highlighted. Results are averaged over 5 runs.

In Table 11, the CECE results further reinforce the effectiveness of CC across all metrics. For CIFAR-10, CC improves CECE for all models compared to baseline methods. For instance, with ResNet-50, the CECE decreases from 0.91 to 0.39. Similar trends are observed on CIFAR-100, with Wide-ResNet showing a reduction in CECE from 0.34 (baseline) to 0.18 when using CC, demonstrating enhanced class-wise calibration.

Dataset	Model	Cross-Entropy		Brier Loss		MMCE		LS-0.05		FLSD-53		FL-3	
		base	ours	base	ours	base	ours	base	ours	base	ours	base	ours
CIFAR-10	ResNet-50	0.91	<b>0.39</b>	0.46	<b>0.35</b>	0.94	<b>0.47</b>	0.71	<b>0.53</b>	0.42	<b>0.35</b>	0.43	<b>0.39</b>
	ResNet-110	0.92	<b>0.41</b>	0.59	<b>0.41</b>	1.04	<b>0.50</b>	<b>0.66</b>	0.67	0.48	<b>0.39</b>	0.43	<b>0.37</b>
	DenseNet-121	0.92	<b>0.43</b>	0.46	<b>0.37</b>	1.04	<b>0.59</b>	0.60	<b>0.48</b>	0.41	<b>0.35</b>	0.42	<b>0.35</b>
	Wide-ResNet	0.68	<b>0.32</b>	0.44	<b>0.32</b>	0.70	<b>0.38</b>	0.79	<b>0.41</b>	0.41	<b>0.28</b>	0.44	<b>0.30</b>
CIFAR-100	ResNet-50	0.38	<b>0.20</b>	0.22	<b>0.19</b>	0.34	<b>0.18</b>	0.23	<b>0.22</b>	0.20	<b>0.19</b>	0.20	<b>0.19</b>
	ResNet-110	0.41	<b>0.21</b>	0.24	<b>0.19</b>	0.42	<b>0.20</b>	0.26	<b>0.22</b>	0.24	<b>0.19</b>	0.24	<b>0.20</b>
	DenseNet-121	0.45	<b>0.23</b>	0.20	<b>0.20</b>	0.42	<b>0.23</b>	0.29	<b>0.22</b>	0.19	0.19	0.20	<b>0.19</b>
	Wide-ResNet	0.34	<b>0.18</b>	0.19	<b>0.18</b>	0.30	<b>0.17</b>	0.21	<b>0.19</b>	0.18	<b>0.17</b>	0.18	<b>0.17</b>

Table 11: Comparison of Train-time Calibration Methods Using CECE↓ Across Various Datasets and Models. CECE values are reported with 15 bins. The best results for each combination is in bold, and our method (CC) is highlighted. Results are averaged over 5 runs.

Table 12 presents the NLL comparison. It is interesting as mentioned in the main section, the CC method sometimes produces higher NLL values.

Dataset	Model	Cross-Entropy		Brier Loss		MMCE		LS-0.05		FLSD-53		FL-3	
		Base	Ours	Base	Ours	Base	Ours	Base	Ours	Base	Ours	Base	Ours
CIFAR-10	ResNet-50	41.2	<b>20.4</b>	<b>18.7</b>	22.3	44.8	<b>20.9</b>	<b>27.7</b>	29.3	<b>17.6</b>	22.7	<b>18.4</b>	24.2
	ResNet-110	47.5	<b>25.5</b>	<b>20.4</b>	22.5	55.7	<b>25.5</b>	29.9	<b>29.4</b>	<b>18.5</b>	21.9	<b>17.8</b>	23.1
	DenseNet-121	42.9	<b>24.0</b>	<b>19.1</b>	21.2	52.1	<b>31.2</b>	28.7	<b>28.5</b>	<b>18.4</b>	27.2	<b>18.0</b>	28.3
	Wide-ResNet	26.8	<b>17.1</b>	<b>15.9</b>	16.2	28.5	<b>18.2</b>	<b>21.7</b>	24.5	<b>14.6</b>	17.6	<b>15.2</b>	19.9
CIFAR-100	ResNet-50	153.7	<b>113.0</b>	<b>99.6</b>	133.5	125.3	<b>116.7</b>	<b>121.0</b>	133.9	<b>88.0</b>	128.8	<b>87.5</b>	128.1
	ResNet-110	179.2	<b>122.3</b>	<b>110.7</b>	146.9	180.6	<b>125.3</b>	<b>133.1</b>	141.4	<b>89.9</b>	126.9	<b>90.9</b>	132.0
	DenseNet-121	205.6	<b>163.1</b>	<b>98.3</b>	139.9	166.6	<b>146.8</b>	<b>142.0</b>	185.8	<b>85.5</b>	129.0	<b>87.1</b>	130.8
	Wide-ResNet	140.1	<b>102.5</b>	<b>84.6</b>	98.7	119.6	<b>109.3</b>	<b>108.1</b>	136.6	<b>76.9</b>	108.7	<b>74.7</b>	106.8

Table 12: Comparison of Train-time Calibration Methods Using NLL↓ Across Various Datasets and Models. The best-performing method for each dataset-model combination is in bold, and our method (CC) is highlighted. Results are averaged over 5 runs.

Table 13 presents a comparison of classification accuracies. While achieving superior calibration performance by CC, the accuracy remains unaffected across all metrics.

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

Dataset	Model	Cross-Entropy		Brier Loss		MMCE		LS-0.05		FLSD-53		FL-3	
		base	ours	base	ours	base	ours	base	ours	base	ours	base	ours
CIFAR-10	ResNet-50	<b>95.05</b>	95.06	<b>94.99</b>	95.01	95.01	<b>94.99</b>	94.71	<b>94.68</b>	95.02	<b>94.95</b>	94.75	94.75
	ResNet-110	<b>95.11</b>	95.16	94.52	<b>94.48</b>	<b>94.60</b>	94.63	<b>94.48</b>	94.49	<b>94.57</b>	94.63	<b>94.92</b>	94.94
	DenseNet-121	95.02	<b>95.01</b>	94.90	<b>94.86</b>	<b>94.59</b>	94.60	94.91	94.91	94.58	<b>94.51</b>	94.66	94.66
	Wide-ResNet	96.13	<b>96.12</b>	95.92	<b>95.90</b>	96.09	<b>96.05</b>	<b>95.80</b>	95.83	<b>95.99</b>	96.01	95.87	95.87
CIFAR-100	ResNet-50	<b>76.70</b>	76.71	76.60	<b>76.58</b>	76.80	76.80	<b>76.56</b>	76.65	76.79	<b>76.73</b>	<b>77.24</b>	77.34
	ResNet-110	77.27	<b>77.17</b>	74.91	<b>74.79</b>	<b>76.93</b>	76.96	<b>76.57</b>	76.64	<b>77.48</b>	77.49	77.08	<b>77.04</b>
	DenseNet-121	<b>75.47</b>	75.49	<b>76.27</b>	76.30	76.03	76.03	<b>75.94</b>	75.96	77.34	77.34	<b>76.76</b>	76.85
	Wide-ResNet	79.29	<b>79.25</b>	79.43	<b>79.29</b>	79.27	<b>79.23</b>	<b>78.83</b>	78.88	<b>79.91</b>	79.92	<b>80.30</b>	80.34

Table 13: Comparison of Train-time Calibration Methods Using Accuracy $\uparrow$  Across Various Datasets and Models. Top-1 Accuracy values are reported. The best results for each combination is in bold, and our method (CC) is highlighted. Results are averaged over 5 runs.