

---

# Two Sides of Mis-Calibration: Identifying Over and Under-Confidence Prediction for Network Calibration (Supplementary Material)

---

Shuang Ao<sup>1</sup>

Stefan Rueger<sup>2</sup>

Advaith Siddharthan<sup>3</sup>

<sup>1,2,3</sup>Knowledge Media Institute., The Open University, Milton Keynes, UK

## A IDENTIFYING UNDER-CONFIDENCE

Table A1: Results of average under-confidence mis-calibration score (UC MCS) and average over-confidence miscalibration score (OC MCS) for baseline, TS and our proposed method cwMCS TS. All results are shown in percentage for clarity. Best results for each row are shown in bold. The value in the bracket shows the percentage of class being under or over-confident.

Dataset	Model	Baseline		TS		cwMCS TS	
		UC MCS (%)	OC MCS (%)	UC MCS (%)	OC MCS (%)	UC MCS (%)	OC MCS (%)
IN	ViT	-4.2 (48.3)	7.4 (51.7)	-5.2 (61.3)	7.2 (38.7)	<b>-3.7 (58.3)</b>	<b>0.5 (41.7)</b>
	SwinT	-11.6 (85.8)	9.2 (14.2)	-5.3 (64)	9.2 (36)	<b>-0.5 (65.7)</b>	<b>8.7 (34.3)</b>
	DeiT	-10.4 (79.6)	9.0 (20.4)	-5.2 (54.7)	9.6 (45.3)	<b>-3.2 (55.3)</b>	<b>5.6 (44.7)</b>
	CaiT	-6.7 (67.6)	9.8 (32.4)	-5.2 (57.6)	9.6 (42.4)	<b>-3.2 (58.2)</b>	<b>7.6 (41.8)</b>
	BeiT	-9.0 (81.9)	8.9 (18.1)	-5.4 (65.3)	8.5 (34.7)	<b>-4.1 (62.2)</b>	<b>8.2 (37.8)</b>
	CoaT	-10.9 (83.1)	8.7 (16.9)	-6.8 (55.6)	9.0 (44.4)	<b>-5.4 (57.8)</b>	<b>8.0 (42.2)</b>
	CrossViT	-9.6 (76.8)	9.5 (23.2)	-5.5 (56.1)	9.6 (43.9)	<b>-3.5 (56.1)</b>	<b>8.6 (43.9)</b>
	ConvMix	-19.0 (90.4)	<b>8.1 (9.6)</b>	-8.9 (61.5)	8.6 (38.5)	<b>-5.7 (59.7)</b>	8.4 (40.3)
	ConvNext	-5.9 (59.8)	9.3 (40.2)	-5.3 (53.8)	9.2 (46.2)	<b>-4.3 (51.6)</b>	<b>9.0 (48.4)</b>
	ResNet34	<b>-4.8 (40.1)</b>	9.9 (59.9)	-6.4 (53.8)	8.6 (46.2)	-6.3 (52)	<b>8.5 (48)</b>
	DenseNet121	<b>-5.2 (43.7)</b>	9.3 (56.3)	-6.3 (55.4)	<b>8.4 (44.6)</b>	-6.4 (55.2)	<b>8.4 (44.8)</b>
Tiny-IN	VGG16	<b>-5.3 (40.5)</b>	8.8 (59.5)	-6.3 (53.9)	8.1 (46.1)	-6.2 (54)	<b>8.0 (46)</b>
	EfficientNet	-17.0 (90.3)	8.4 (9.7)	-16.0 (87.4)	8.2 (12.6)	<b>-8.8 (53.6)</b>	<b>1.2 (46.4)</b>
	ResNet34	<b>-3.1 (11.5)</b>	10.6 (88.5)	-5.8 (58)	6.3 (42)	-5.4 (52)	<b>4.3 (48)</b>
	DenseNet121	-1.0 (2.5)	13.4 (97.5)	-5.5 (58.5)	7.6 (41.5)	<b>-5.2 (56.5)</b>	<b>6.6 (43.5)</b>
C100	VGG16	<b>-0.8 (1.5)</b>	15.9 (98.5)	-5.2 (57.5)	5.1 (42.5)	-4.9 (52.5)	<b>3.1 (47.5)</b>
	Res34	<b>-0.1 (5)</b>	13.4 (95)	-5.2 (61)	6.2 (39)	-5.0 (57)	<b>4.2 (43)</b>
	DenseNet121	<b>-1.3 (1)</b>	15 (99)	-5.7 (55)	6.3 (46)	-4.3 (58)	<b>5.6 (42)</b>
C10	VGG16	<b>-1.3 (3)</b>	10.3 (97)	-4.4 (56)	8.7 (44)	-4.4 (56)	<b>3.7 (44)</b>
	ResNet34	0.0 (0)	4.1 (100)	<b>-1.4 (30)</b>	2.1 (70)	-1.5 (30)	<b>1.5 (70)</b>
	DenseNet121	0.0 (0)	11.3 (100)	-1.6 (40)	5.4 (60)	<b>-0.2 (30)</b>	<b>5.2 (70)</b>
	VGG16	0.0 (0)	7.1 (100)	-0.9 (40)	2.8 (60)	<b>-0.7 (50)</b>	<b>0.2 (50)</b>

Table A1 illustrates the results of mean under-confidence and mean over-confidence scores, as well as the percentage of classes with different confidence statuses correspondingly. For ImageNet dataset with transformers variants, most of the classes are under-confident with baselines, where the absolute value of mean under-confidence score is higher than the mean over-confidence score. The model with the highest percentage of under-confident classes is ConvMix, where only 10 percent of classes are over-confident. When it comes to CNNs, over and under-confident classes are more balanced. Surprisingly, EfficientNet has a similar behavior as ConvMix, where the percentage of under-confident classes are much higher than over-confident ones. After applying TS, the percentage of over and under-confident classes are more balanced, and our proposed method cwMCS TS keeps this trend. Compared to baseline, our cwMCS TS method almost halves over and under-confidence scores, whereas TS only makes a slight change of them. For Tiny-ImageNet, CIFAR100 and CIFAR10

datasets with CNNs, more than 90 percent of classes are over-confident in baselines, with none of the classes under-confident for CIFAR10 dataset. However, more than half of the classes become under-confident after applying TS, indicating that TS can overly calibrate models. Our proposed method cwMCS TS significantly improves the mean over-confidence score and contributes to better calibration for under-confident classes.