

# A Stitch in Time Saves Nine: A Train-Time Regularizing Loss for Improved Neural Network Calibration

Ramya Hebbalaguppe<sup>1,2,\*</sup> Jatin Prakash<sup>1,\*</sup> Neelabh Madan<sup>1,\*</sup> Chetan Arora

<sup>1</sup>Indian Institute of Technology Delhi, India

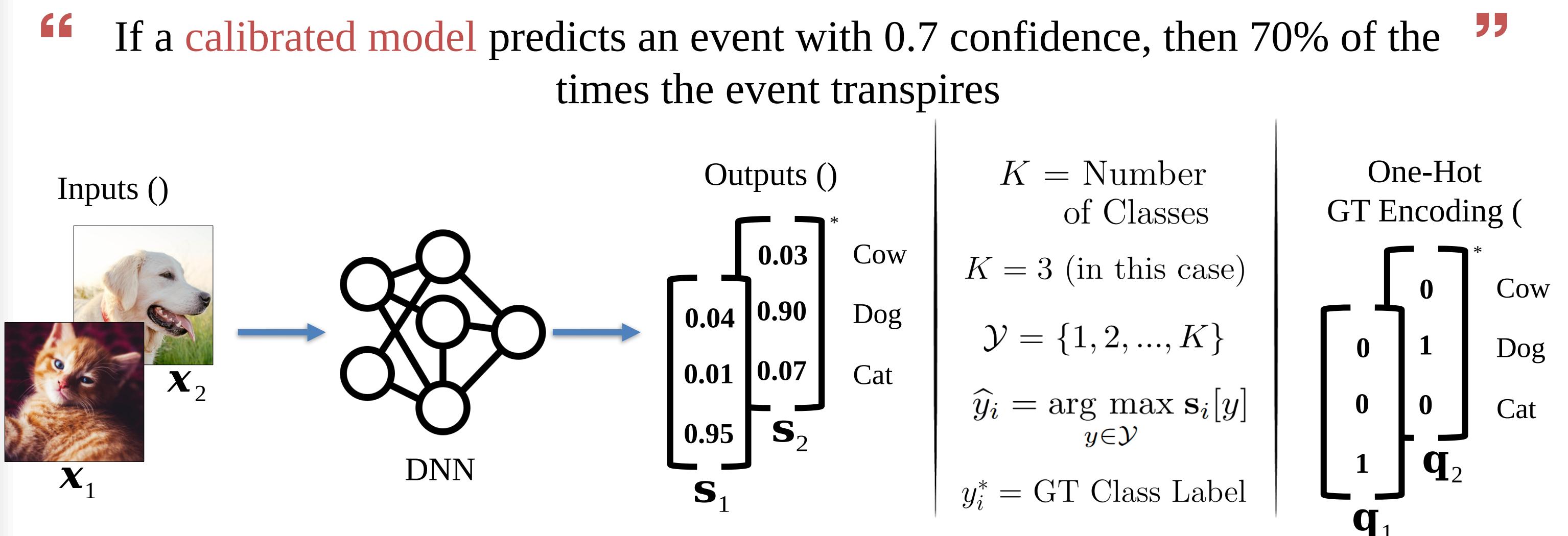
<sup>2</sup>TCS Research, India



## Highlights

- [Novelty] We propose an auxiliary loss to overcome miscalibration
- [Multi-class Calibration] Our method takes into account the entire probability vector
- [Powerful Regularizer] Models trained using our method are well calibrated even under domain/dataset drift
- [Superior Calibration] Outperforms SOTA methods on various datasets and models
- [Beyond Image Classification] Promising results in semantic segmentation in images and NL classification tasks

## Understanding Calibration



### Top-Label Calibration

$$\mathbb{P}(\hat{y}_i = y_i^* \mid \mathbf{s}[\hat{y}_i] = p) = p$$

### Multi-class Calibration

$$\mathbb{P}(y = y_i^* \mid \mathbf{s}_i[y] = p) = p \quad \forall y \in \mathcal{Y}$$

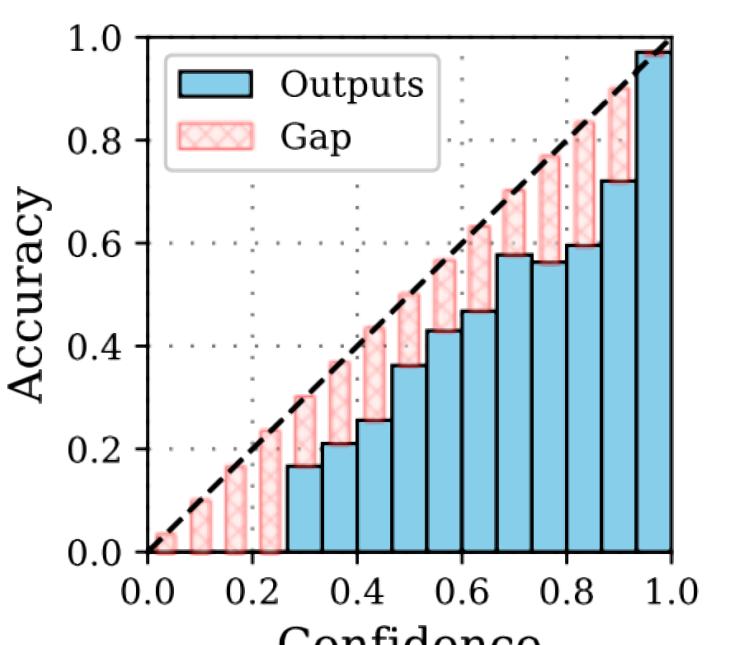
**Problem:** Modern Neural Networks are neither top-label nor multi-class calibrated

## Measuring Calibration

### 1. Quantitative Measures

- [ECE] Expected Calibration Error: It calculates the absolute difference between the model's accuracy and confidence. It captures the information about top-label calibration.
- [SCE] Static Calibration Error: A simple class-wise extension to ECE that captures multi-class calibration

### 2. Reliability Diagrams



Paper and Code: [github.com/mdca-loss](https://github.com/mdca-loss)

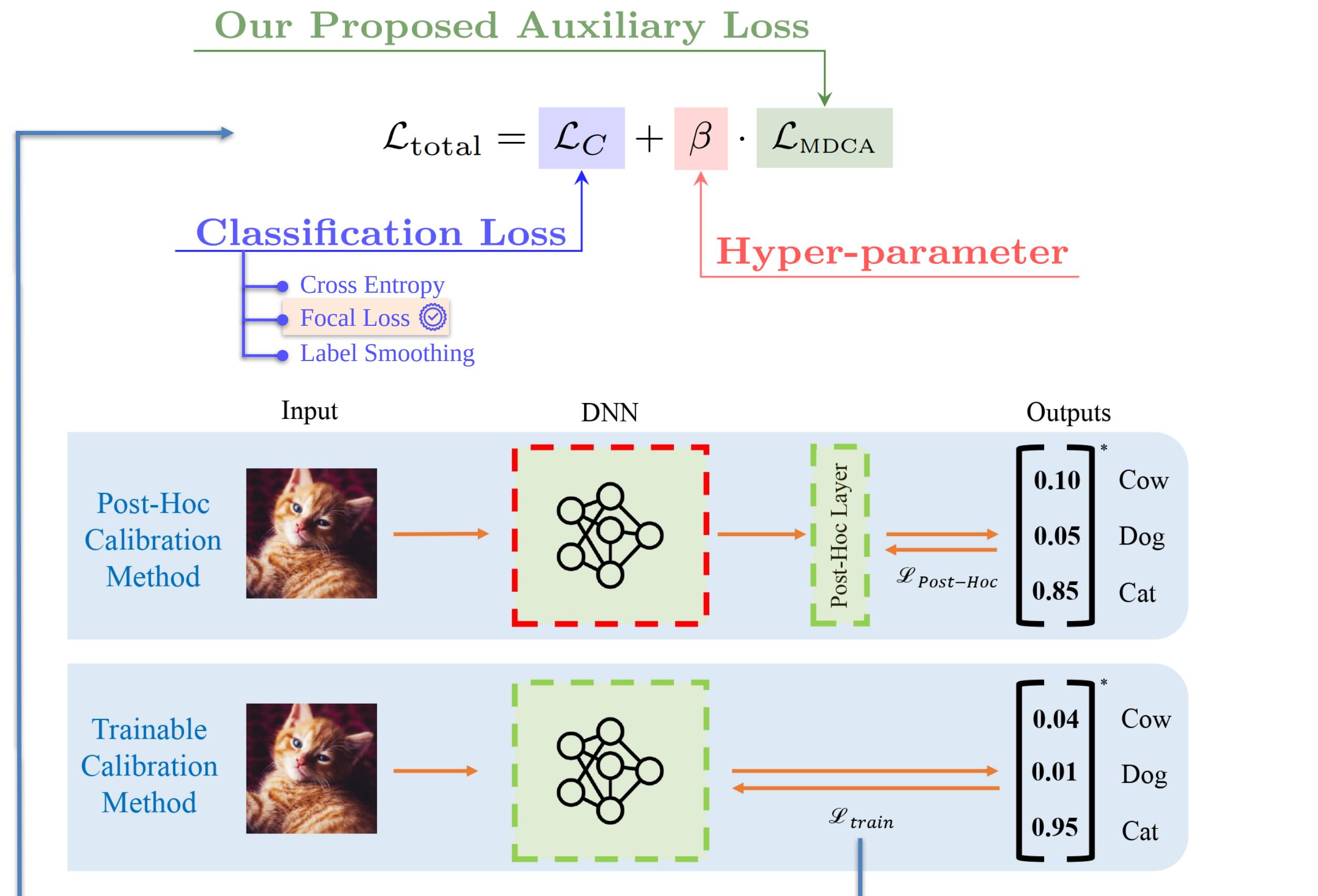


## Proposed Solution

We propose a novel train-time regularizing auxiliary loss function called **Multi-class Difference in Confidence and Accuracy (MDCA)**

$$\mathcal{L}_{\text{MDCA}} = \frac{1}{K} \sum_{j=1}^K \left| \frac{1}{N_b} \sum_{i=1}^{N_b} \mathbf{s}_i[j] - \frac{1}{N_b} \sum_{i=1}^{N_b} q_i[j] \right|$$

Avg. Confidence      Avg. Count  
Number of classes      Number of samples in a mini-batch



- Frozen Architecture
  - Trainable Architecture
  - Post-Hoc Calibration Method
  - Trainable Calibration Method
- Hold out set required
  - Only a few parameters available  $O(1)/O(K)/O(K^2)$
  - No hold-out set required
  - Millions/Billions of parameters available for calibration
- $K = \text{Number of Classes}$

\* The output confidence values are for illustration purposes only

## Experimental Results

### 1. Superior performance against trainable calibration methods

Dataset	Model	BS [2]			DCA [31]			MMCE [26]			FLSD [37]			Ours (FL+MDCA)		
		SCE	ECE	TE	SCE	ECE	TE	SCE	ECE	TE	SCE	ECE	TE	SCE	ECE	TE
CIFAR10	ResNet32	6.60	2.92	7.76	8.41	4.00	<b>7.06</b>	8.17	3.31	8.41	9.48	4.41	7.87	<b>3.22</b>	<b>0.93</b>	7.18
	ResNet56	5.44	2.17	7.75	7.59	3.38	<b>6.53</b>	9.11	3.71	8.23	7.71	3.49	7.04	<b>2.93</b>	<b>0.70</b>	7.08
CIFAR100	ResNet32	1.97	5.32	33.53	2.82	11.31	29.67	2.79	11.09	31.62	1.77	1.69	32.15	<b>1.72</b>	<b>1.49</b>	<b>31.58</b>
	ResNet56	1.86	4.69	30.72	2.77	9.29	43.43	2.35	8.61	28.75	1.71	1.90	29.11	<b>1.60</b>	<b>0.23</b>	3.85
SVHN	ResNet20	2.12	<b>0.45</b>	<b>3.56</b>	4.29	2.02	3.83	9.18	4.34	4.12	18.98	9.37	4.10	<b>1.90</b>	0.47	3.92
	ResNet56	2.18	0.66	<b>3.25</b>	2.16	0.49	3.32	9.69	4.48	4.26	26.15	13.23	3.65	<b>1.51</b>	<b>0.23</b>	3.85
Mendeley V2	ResNet50	117.6	3.75	18.43	145.1	8.29	17.47	130.4	<b>3.45</b>	<b>15.06</b>	104.3	9.64	19.71	<b>85.68</b>	4.81	17.95
Tiny-ImageNet	ResNet34	1.53	7.79	43.00	2.11	17.40	<b>36.68</b>	1.62	9.71	40.75	1.18	<b>1.91</b>	37.01	<b>1.17</b>	1.99	37.49
20 Newsgroups	Global-Pool CNN	725.82	13.71	<b>25.93</b>	719.83	15.30	28.07	731.31	12.69	28.63	940.70	<b>4.52</b>	30.80	<b>487.82</b>	16.55	27.88

### 2. Superior class-wise calibration

Method	Classes									
	0	1	2	3	4	5	6	7	8	
Cross Entropy	<b>0.20</b>	0.62	0.33	0.65	0.23	0.36	0.25	0.26	<b>0.21</b>	0.41
Focal Loss [32]	0.30	0.48	0.41	<b>0.18</b>	0.38	0.19	0.33	0.36	0.32	0.30
LS [38]	1.63	2.60	2.54	1.90	1.91	1.74	1.73	1.75	1.63	1.58
Brier Score [2]	0.23	0.28	0.40	0.45	0.25	0.26	0.25	0.27	0.21	0.37
MMCE [26]	1.78	2.35	2.12	2.00	1.74	1.87	1.65	1.76	1.70	1.84
DCA [31]	0.31	0.70	0.40	0.72	0.31	0.46	0.35	0.35	0.37	0.36
FLSD [37]	1.52	3.24	2.74	2.15	1.79	1.82	1.84	1.62	1.54	1.38
Ours (FL+MDCA)	0.22	<b>0.16</b>	<b>0.24</b>	0.25	<b>0.22</b>	<b>0.16</b>	<b>0.16</b>	<b>0.17</b>	0.25	<b>0.20</b>

### 3. Performance under dataset drift

Method	Art	Cartoon	Sketch	Average
NLL	6.33	17.95	15.01	13.10
LS [38]	7.80	11.95	<b>10.88</b>	10.21
FL	8.61	16.62	10.94	12.06
Brier Score [2]	6.55	13.19	15.63	11.79
MMCE [26]	6.35	15.70	17.16	13.07
DCA [31]	7.49	18.01	14.99	13.49
FLSD [37]	8.35	13.39	13.86	11.87
Ours (FL+MDCA)	<b>6.21</b>	<b>11.91</b>	11.08	<b>9.73</b>

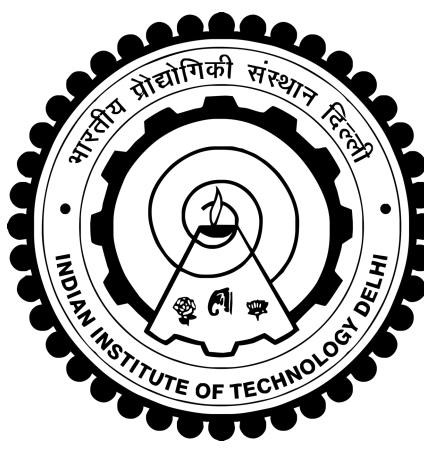
### 4. Mitigating overconfident mistakes



Sampled incorrect predictions only

### 5. Performance under data imbalance

Method	IF-10	CIFAR10 IF-50	IF-100	SVHN IF-2.7
NLL	18.44	32.21	31.04	3.43
FL [32]	14.65	29.67	28.89	2.54
LS [38]	14.88	26.30	<b>20.79</b>	18.80
BS [2]	15.74	33.57	29.01	2.12
MMCE [26]	15.10	29.05	21.56	9.18
FLSD [37]	16.05	31.35	30.28	18.98
DCA [31				



# A Stitch in Time Saves Nine: A Train-Time Regularizing Loss for Improved Neural Network Calibration

Ramya Hebbalaguppe<sup>1,2,\*</sup> Jatin Prakash<sup>1,\*</sup> Neelabh Madan<sup>1,\*</sup> Chetan Arora

<sup>1</sup>Indian Institute of Technology Delhi, India

<sup>2</sup>TCS Research, India

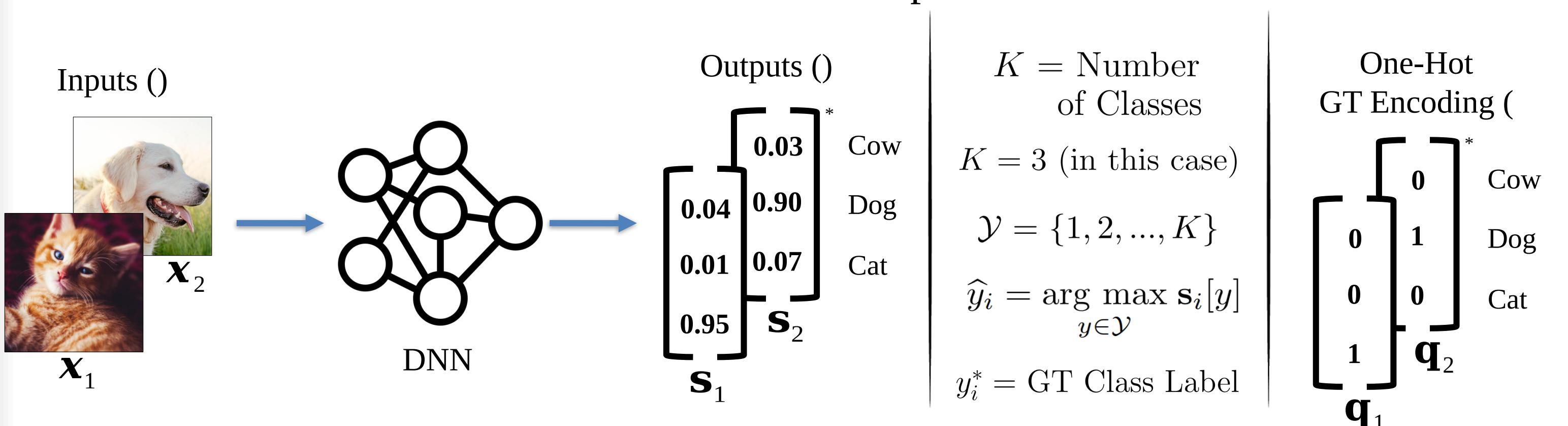


## Highlights

- [Novelty] We propose an auxiliary loss to overcome miscalibration
- [Multi-class Calibration] The entire probability vector (all K classes) taken into account
- [Powerful Regularizer] Models trained using our method are relatively well calibrated even under domain/dataset drift
- [Superior Calibration] Outperforms SOTA methods on various datasets and models
- [Beyond Image Classification] Promising results in semantic segmentation in images and NL classification tasks

## Understanding Calibration

If a calibrated model predicts an event with 0.7 confidence, then 70% of the times the event transpires



### Top-Label Calibration

$$\mathbb{P}(\hat{y}_i = y_i^* \mid \mathbf{s}[\hat{y}_i] = p) = p$$

$$\mathbb{P}(y = y_i^* \mid \mathbf{s}_i[y] = p) = p \quad \forall y \in \mathcal{Y}$$

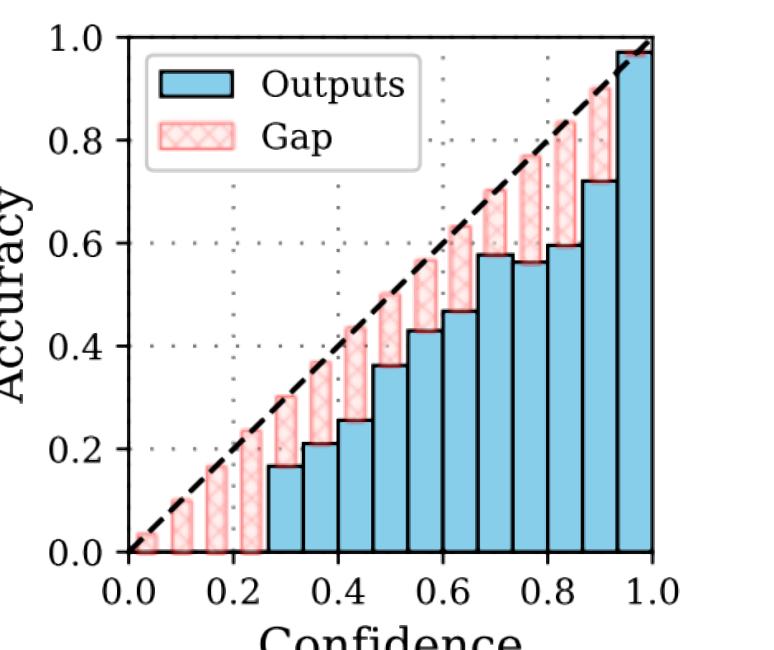
**Problem:** Modern Neural Networks are neither top-label nor multi-class calibrated

## Measuring Calibration

### 1. Quantitative Measures

- [ECE] Expected Calibration Error: It calculates the absolute difference between the model's accuracy and confidence. It captures the information about top-label calibration.
- [SCE] Static Calibration Error: A simple class-wise extension to ECE that captures multi-class calibration

### 2. Reliability Diagrams



Paper and Code: [github.com/mdca-loss](https://github.com/mdca-loss)



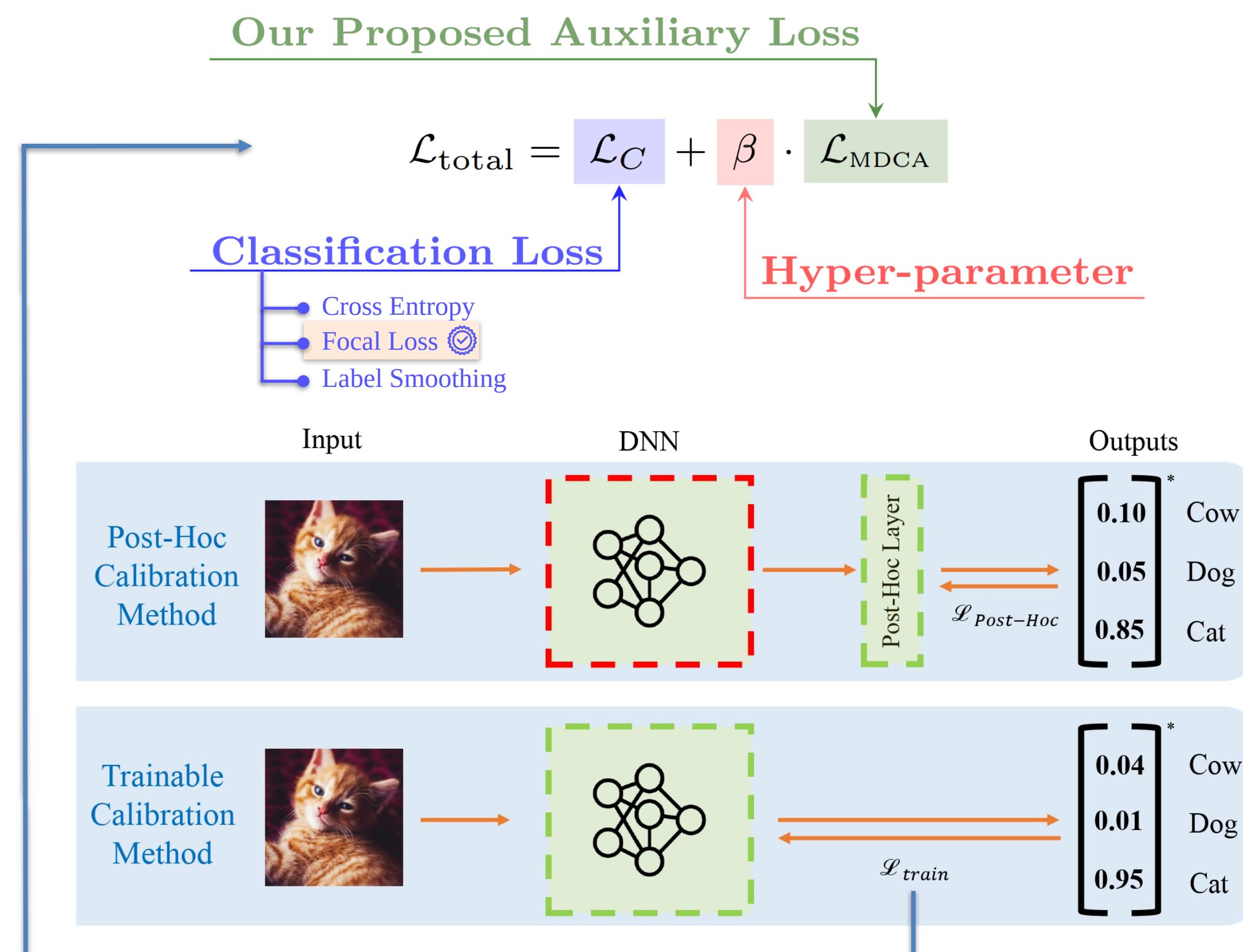
## Proposed Solution

We propose a novel train-time regularizing auxiliary loss function called **Multi-class Difference in Confidence and Accuracy (MDCA)**

$$\mathcal{L}_{\text{MDCA}} = \frac{1}{K} \sum_{j=1}^K \left| \frac{1}{N_b} \sum_{i=1}^{N_b} \mathbf{s}_i[j] - \frac{1}{N_b} \sum_{i=1}^{N_b} q_i[j] \right|$$

Avg. Confidence      Avg. Count

Number of classes      Number of samples in a mini-batch



- |                        |                             |                              |
|------------------------|-----------------------------|------------------------------|
| Frozen Architecture    | Post-Hoc Calibration Method | Trainable Calibration Method |
| Dashed Box             | Post-Hoc Calibration Method | Trainable Calibration Method |
| Trainable Architecture | Dashed Box                  | No hold-out set required     |
- Hold out set required
  - Only a few parameters available  $O(1)/O(K)/O(K^2)$
  - No hold-out set required
  - Millions/Billions of parameters available for calibration
- $K = \text{Number of Classes}$

\* The output confidence values are for illustration purposes only

## Experimental Results

### 1. Superior performance against trainable calibration methods

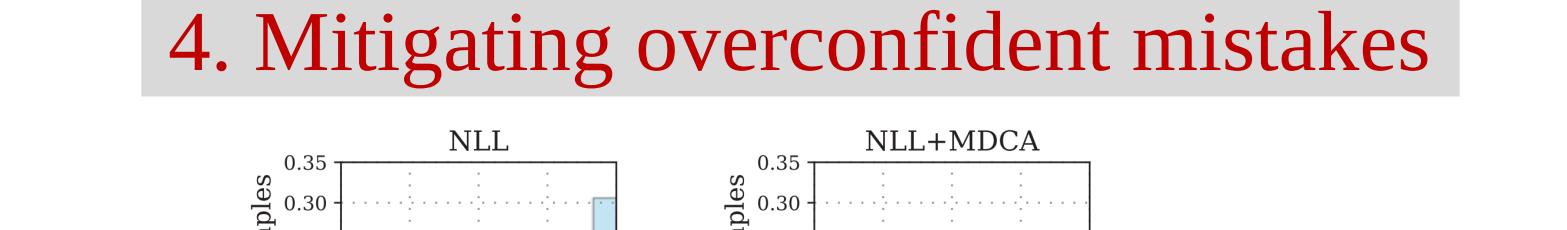
Dataset	Model	BS [2]			DCA [31]			MMCE [26]			FLSD [37]			Ours (FL+MDCA)		
		SCE	ECE	TE	SCE	ECE	TE	SCE	ECE	TE	SCE	ECE	TE	SCE	ECE	TE
CIFAR10	ResNet32	6.60	2.92	7.76	8.41	4.00	<b>7.06</b>	8.17	3.31	8.41	9.48	4.41	7.87	<b>3.22</b>	<b>0.93</b>	7.18
	ResNet56	5.44	2.17	7.75	7.59	3.38	<b>6.53</b>	9.11	3.71	8.23	7.71	3.49	7.04	<b>2.93</b>	<b>0.70</b>	7.08
CIFAR100	ResNet32	1.97	5.32	33.53	2.82	11.31	29.67	2.79	11.09	31.62	1.77	1.69	32.15	<b>1.72</b>	<b>1.49</b>	<b>31.58</b>
	ResNet56	1.86	4.69	30.72	2.77	9.29	43.43	2.35	8.61	28.75	1.71	1.90	29.11	<b>1.60</b>	<b>0.23</b>	3.85
SVHN	ResNet20	2.12	<b>0.45</b>	<b>3.56</b>	4.29	2.02	3.83	9.18	4.34	4.12	18.98	9.37	4.10	<b>1.90</b>	0.47	3.92
	ResNet56	2.18	0.66	<b>3.25</b>	2.16	0.49	3.32	9.69	4.48	4.26	26.15	13.23	3.65	<b>1.51</b>	<b>0.23</b>	3.85
Mendeley V2	ResNet50	117.6	3.75	18.43	145.1	8.29	17.47	130.4	<b>3.45</b>	<b>15.06</b>	104.3	9.64	19.71	<b>85.68</b>	4.81	17.95
Tiny-ImageNet	ResNet34	1.53	7.79	43.00	2.11	17.40	<b>36.68</b>	1.62	9.71	40.75	1.18	<b>1.91</b>	37.01	<b>1.17</b>	1.99	37.49
20 Newsgroups	Global-Pool CNN	725.82	13.71	<b>25.93</b>	719.83	15.30	28.07	731.31	12.69	28.63	940.70	<b>4.52</b>	30.80	<b>487.82</b>	16.55	27.88

### 2. Superior class-wise calibration

Method	Classes									
	0	1	2	3	4	5	6	7	8	
Cross Entropy	<b>0.20</b>	0.62	0.33	0.65	0.23	0.36	0.25	0.26	<b>0.21</b>	0.41
Focal Loss [32]	0.30	0.48	0.41	<b>0.18</b>	0.38	0.19	0.33	0.36	0.32	0.30
LS [38]	1.63	2.60	2.54	1.90	1.91	1.74	1.73	1.75	1.63	1.58
Brier Score [2]	0.23	0.28	0.40	0.45	0.25	0.26	0.25	0.27	0.21	0.37
MMCE [26]	1.78	2.35	2.12	2.00	1.74	1.87	1.65	1.76	1.70	1.84
DCA [31]	0.31	0.70	0.40	0.72	0.31	0.46	0.35	0.35	0.37	0.36
FLSD [37]	1.52	3.24	2.74	2.15	1.79	1.82	1.84	1.62	1.54	1.38
Ours (FL+MDCA)	0.22	<b>0.16</b>	<b>0.24</b>	0.25	<b>0.22</b>	<b>0.16</b>	<b>0.16</b>	<b>0.17</b>	0.25	<b>0.20</b>

### 3. Performance under dataset drift

Method	Art	Cartoon	Sketch	Average
NLL	6.33	17.95	15.01	13.10
LS [38]	7.80	11.95	<b>10.88</b>	10.21
FL	8.61	16.62	10.94	12.06
Brier Score [2]	6.55	13.19	15.63	11.79
MMCE [26]	6.35	15.70	17.16	13.07
DCA [31]	7.49	18.01	14.99	13.49
FLSD [37]	8.35	13.39	13.86	11.87
Ours (FL+MDCA)	<b>6.21</b>	<b>11.91</b>	11.08	<b>9.73</b>



### 4. Mitigating overconfident mistakes



Sampled incorrect predictions only

### 5. Performance under data imbalance

Method	IF-10	CIFAR10 IF-50	IF-100	SVHN IF-2.7
NLL	18.44	32.21	31.04	3.43
FL [32]	14.65	29.67	28.89	2.54
LS [38]	14.88	26.30	<b>20.79</b>	18.80
BS [2]	15.74	33.57	29.01	2.12
MMCE [26]				