

# Beyond temperature scaling: Obtaining well-calibrated multiclass probabilities with Dirichlet calibration

Meelis Kull, Miquel Perello Nieto, Markus Kängsepp, Telmo Silva  
Filho, Hao Song, Peter Flach

NeurIPS 2019



UNIVERSITY OF TARTU



University of  
BRISTOL



Universidade Federal da Paraíba  
ESTATÍSTICA



The  
Alan Turing  
Institute



# Contributions

- New parametric calibration method:

	Logit space	Class probability space
Binary classification	Derived from Gaussian distribution <b>Platt scaling</b> <sup>[1]</sup>	Derived from Beta distribution <b>Beta calibration</b> <sup>[2]</sup> (+ constrained variants)
Multi-class classification	Matrix scaling <sup>[3]</sup> (+ vector scaling, temperature scaling)	Derived from Dirichlet distribution <b>Dirichlet calibration</b> (+ constrained variants)

- New regularization method for matrix scaling (and for Dirichlet calibration):

ODIR – Off-Diagonal and Intercept Regularisation

- Multi-class classifier evaluation:

Confidence-calibrated  
Classwise-calibrated  
Multiclass-calibrated

Confidence-reliability diagram  
Classwise-reliability diagrams

Confidence-ECE  
Classwise-ECE

# Making classifiers more trustworthy



UNIVERSITY OF TARTU



University of  
BRISTOL



Universidade Federal da Paraíba  
ESTATÍSTICA



The  
Alan Turing  
Institute



# Making classifiers more trustworthy

a classifier with 60% accuracy

on a set of instances

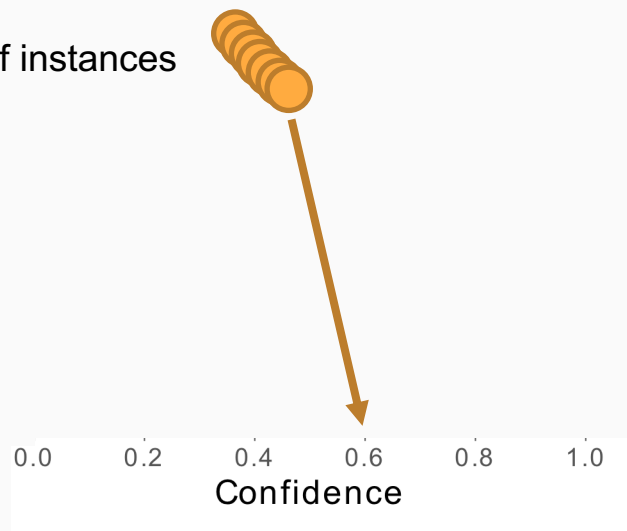




# Making classifiers more trustworthy

a classifier with 60% accuracy

on a set of instances



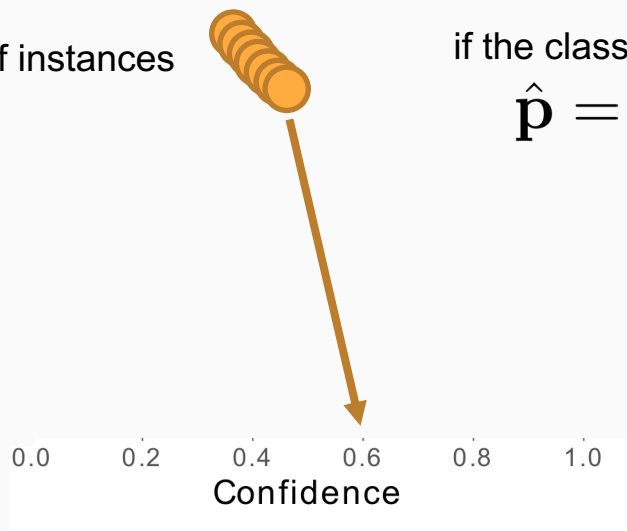
# Making classifiers more trustworthy

a classifier with 60% accuracy

on a set of instances

if the classifier reports class probabilities

$$\hat{\mathbf{p}} = (\hat{p}_1, \hat{p}_2, \dots, \hat{p}_k)$$



# Making classifiers more trustworthy

a classifier with 60% accuracy

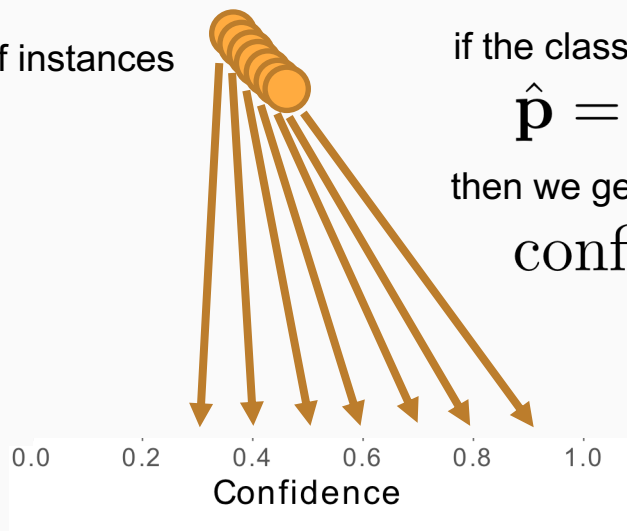
on a set of instances

if the classifier reports class probabilities

$$\hat{\mathbf{p}} = (\hat{p}_1, \hat{p}_2, \dots, \hat{p}_k)$$

then we get instance-specific

$$\text{confidence} = \max(\hat{\mathbf{p}})$$



# Trustworthy if confidence-calibrated

0.0 0.2 0.4 0.6 0.8 1.0  
Confidence



UNIVERSITY OF TARTU



University of  
BRISTOL



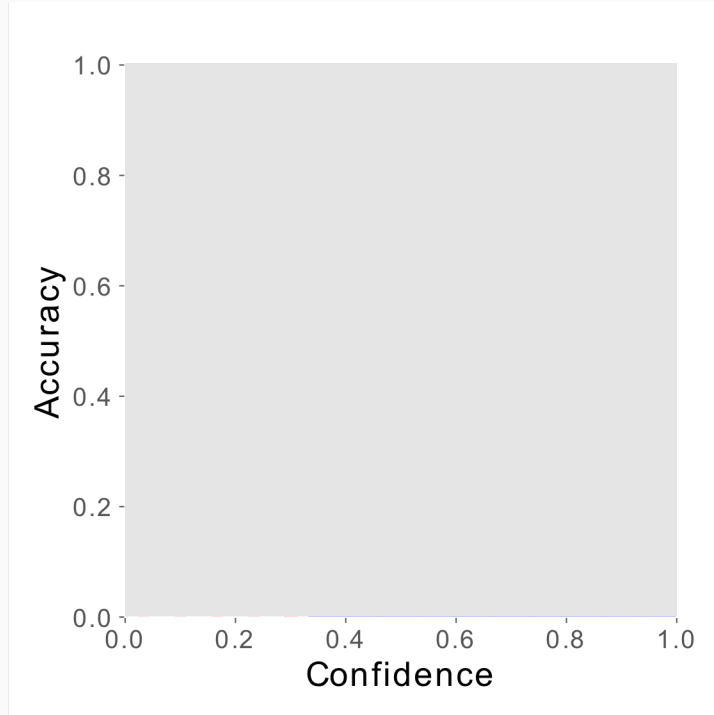
Universidade Federal da Paraíba  
ESTATÍSTICA



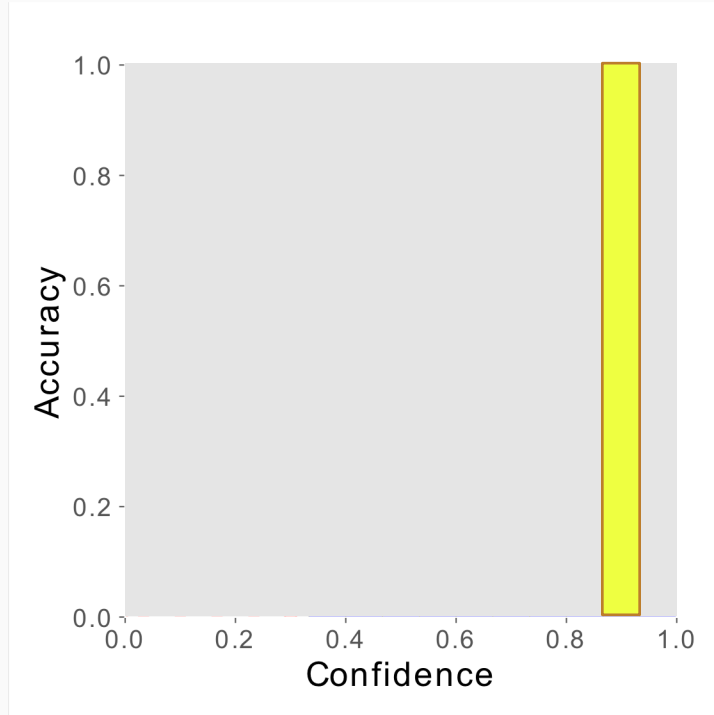
The  
Alan Turing  
Institute



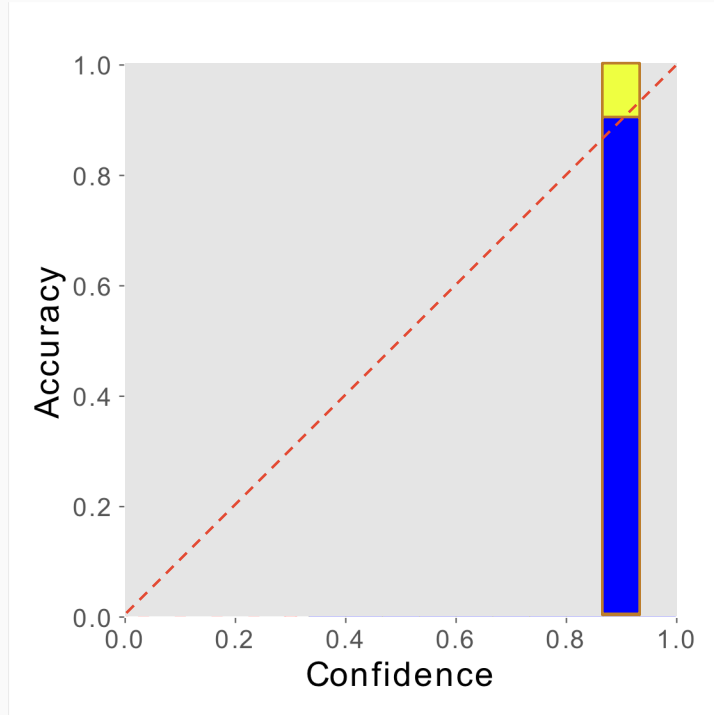
# Trustworthy if confidence-calibrated



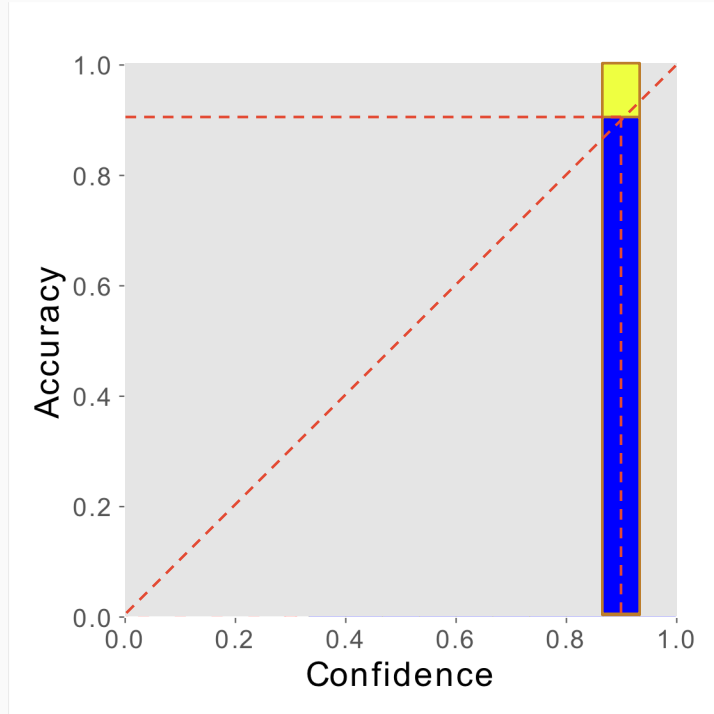
# Trustworthy if confidence-calibrated



# Trustworthy if confidence-calibrated



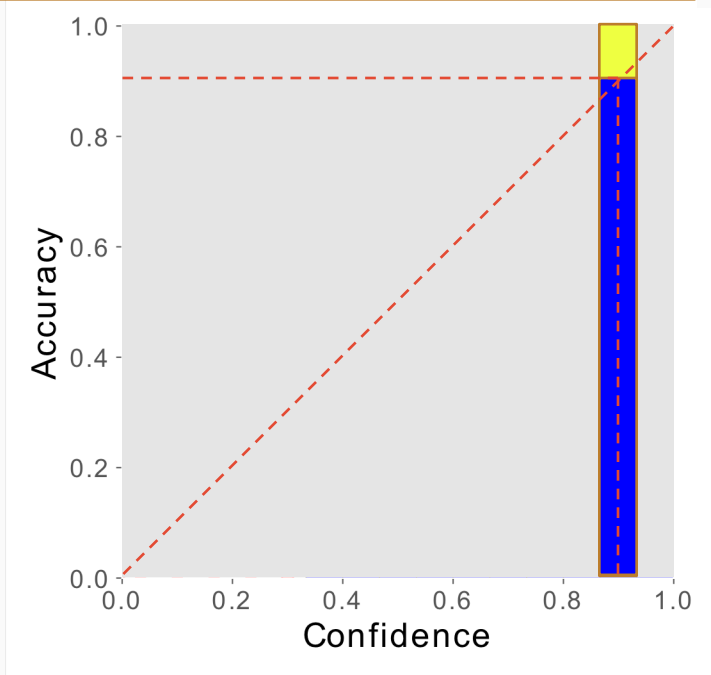
# Trustworthy if confidence-calibrated





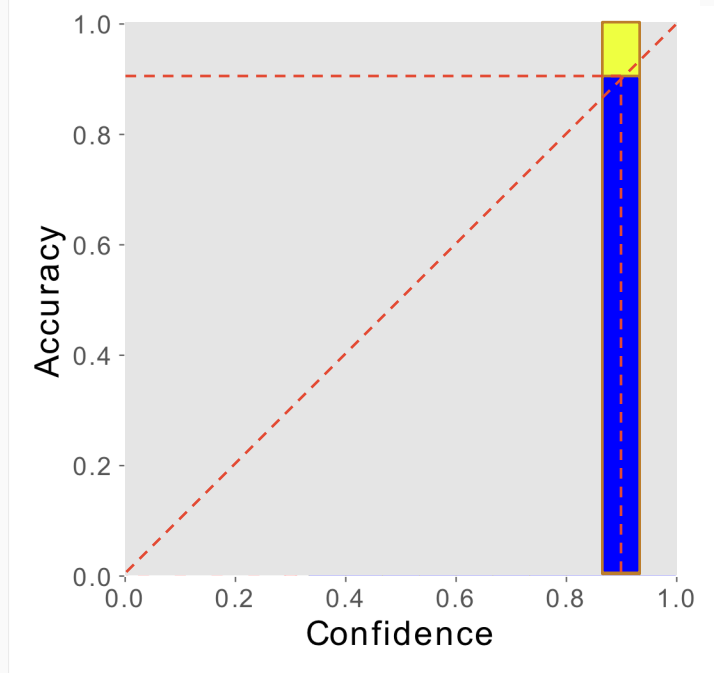
# Trustworthy if confidence-calibrated

$$\max \hat{p}(X) = 0.9$$



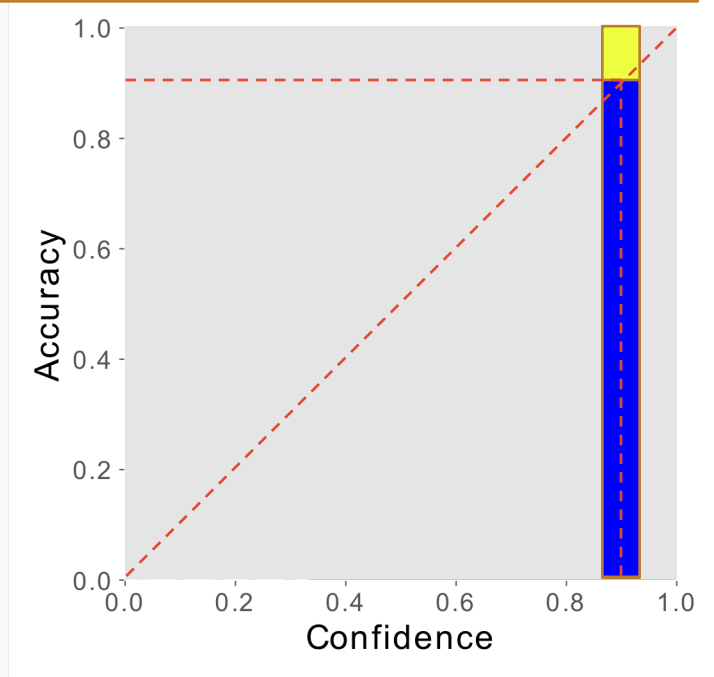
# Trustworthy if confidence-calibrated

$$Y = \arg \max \hat{\mathbf{p}}(X) \quad \max \hat{\mathbf{p}}(X) = 0.9$$



# Trustworthy if confidence-calibrated

$$P(Y = \arg \max \hat{\mathbf{p}}(X) \mid \max \hat{\mathbf{p}}(X) = 0.9) = 0.9$$

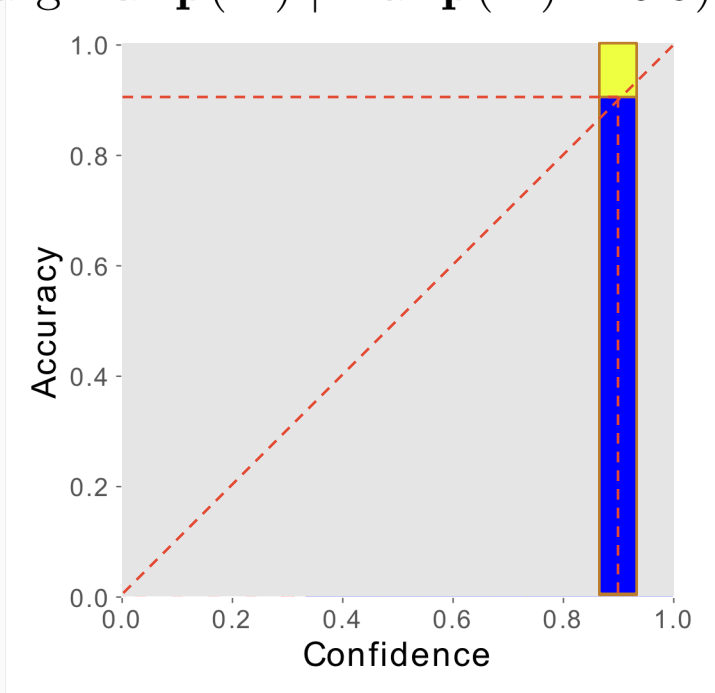


# Trustworthy if confidence-calibrated

$$P(Y = \arg \max \hat{\mathbf{p}}(X) \mid \max \hat{\mathbf{p}}(X) = 0.9) = 0.9$$

Confidence-calibrated:

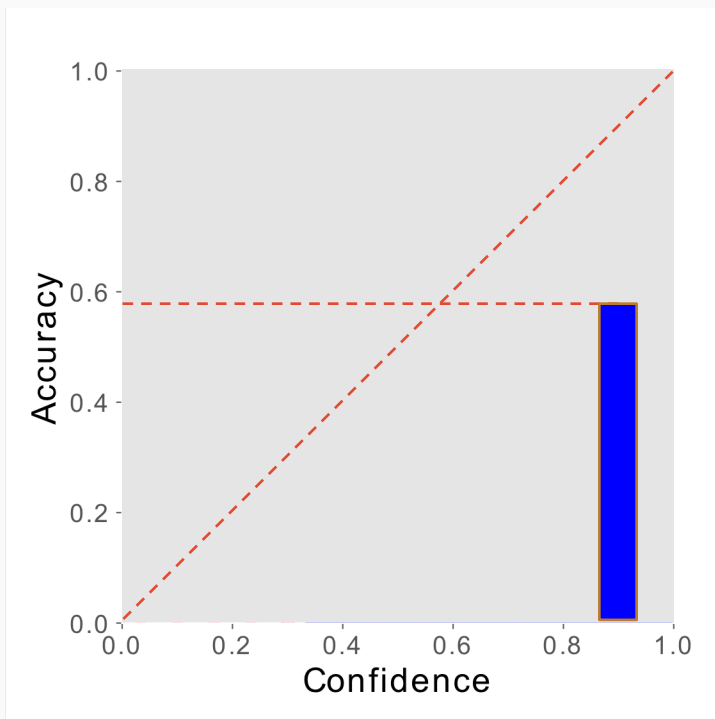
$$P(Y = \arg \max \hat{\mathbf{p}}(X) \mid \max \hat{\mathbf{p}}(X) = c) = c$$



# Deep nets are usually over-confident

Confidence-calibrated:

$$P(Y = \arg \max \hat{\mathbf{p}}(X) \mid \max \hat{\mathbf{p}}(X) = c) = c$$



Experimental setup:

CIFAR-10

ResNet Wide 32

Accuracy:

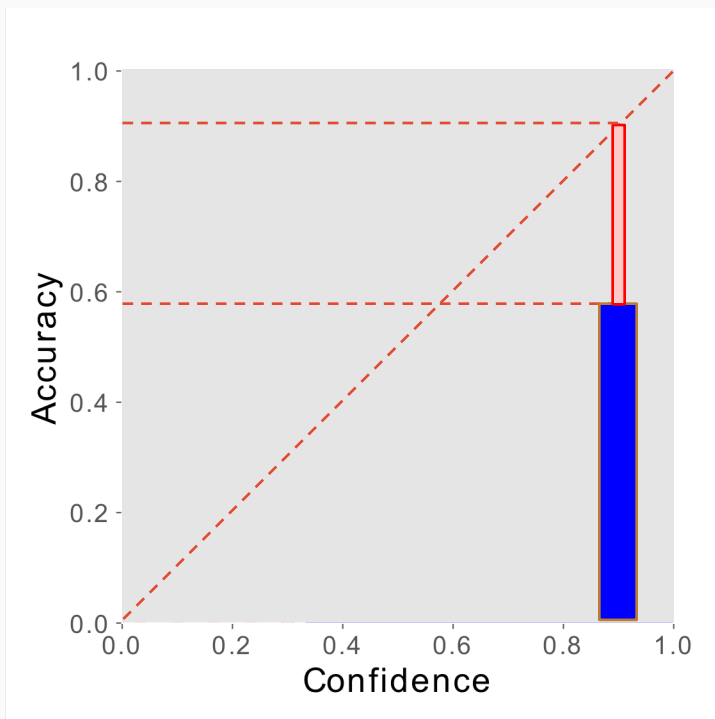
Overall: 94%

At 90% confidence: 58%

# Deep nets are usually over-confident

Confidence-calibrated:

$$P(Y = \arg \max \hat{\mathbf{p}}(X) \mid \max \hat{\mathbf{p}}(X) = c) = c$$



Experimental setup:

CIFAR-10

ResNet Wide 32

Accuracy:

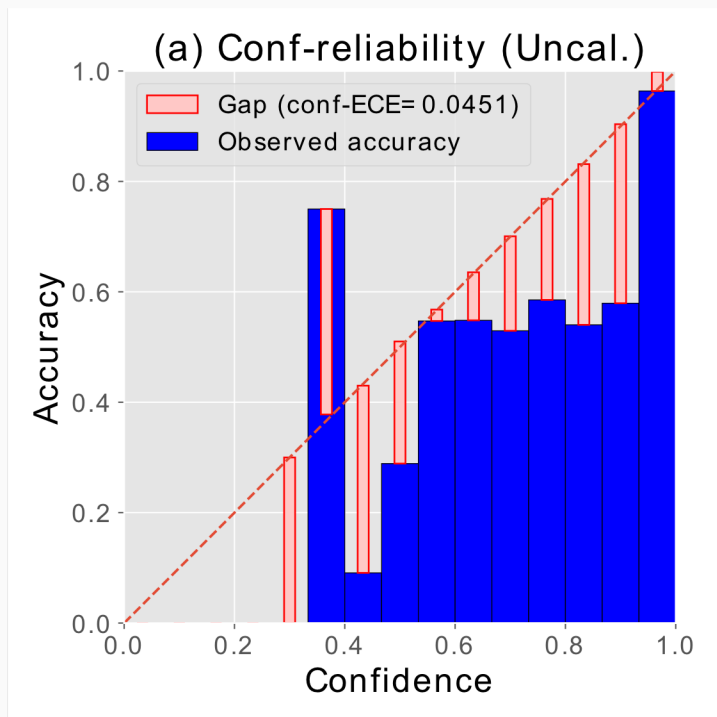
Overall: 94%

At 90% confidence: 58%

# Example: uncalibrated predictions

Confidence-calibrated:

$$P(Y = \arg \max \hat{\mathbf{p}}(X) \mid \max \hat{\mathbf{p}}(X) = c) = c$$



Experimental setup:

CIFAR-10

ResNet Wide 32

Accuracy:

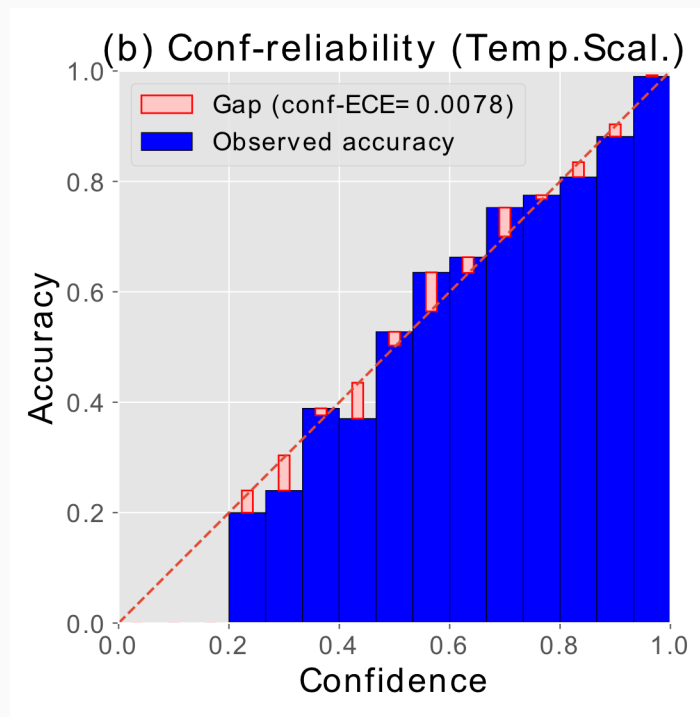
Overall: 94%

At 90% confidence: 58%

# Example: after calibration with temperature scaling

Confidence-calibrated:

$$P(Y = \arg \max \hat{\mathbf{p}}(X) \mid \max \hat{\mathbf{p}}(X) = c) = c$$



Experimental setup:

CIFAR-10

ResNet Wide 32

Accuracy:

Overall: 94%

At 90% confidence: 58%

Accuracy after Temp.Scal:

Overall: 94%

At 90% confidence: 88%



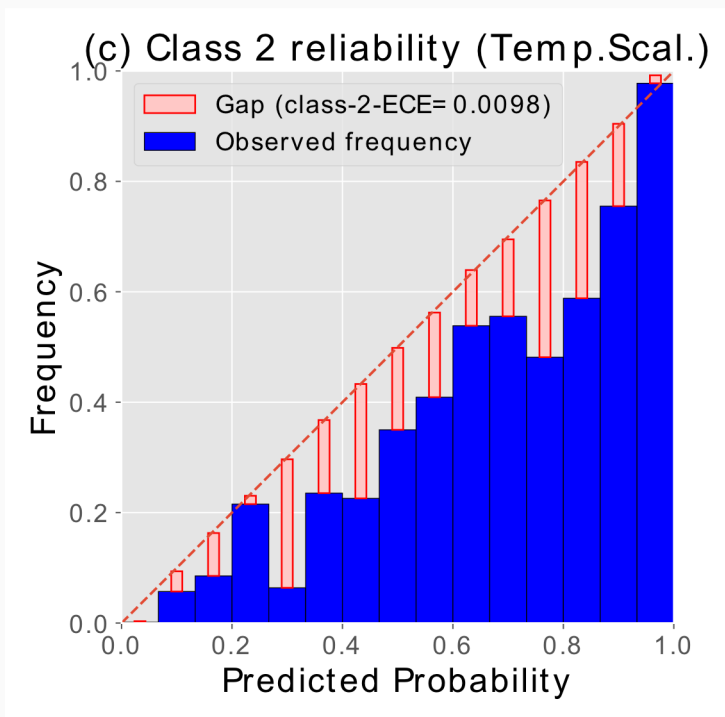
# Example: after calibration with temperature scaling

Confidence-calibrated:

$$P(Y = \arg \max \hat{\mathbf{p}}(X) \mid \max \hat{\mathbf{p}}(X) = c) = c$$

Classwise-calibrated:

$$P(Y = i \mid \hat{p}_i(X) = c) = c$$



Experimental setup:

CIFAR-10

ResNet Wide 32

Accuracy:

Overall: 94%

At 90% confidence: 58%

Accuracy after Temp.Scal:

Overall: 94%

At 90% confidence: 88%

At 90% class 2 prob: 76%

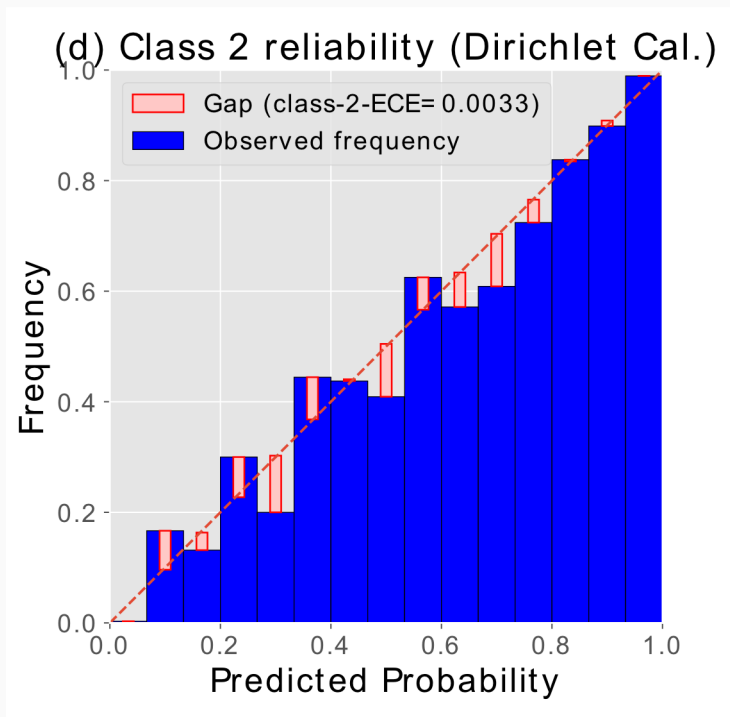
# Example: after calibration with Dirichlet calibration

Confidence-calibrated:

$$P(Y = \arg \max \hat{\mathbf{p}}(X) \mid \max \hat{\mathbf{p}}(X) = c) = c$$

Classwise-calibrated:

$$P(Y = i \mid \hat{p}_i(X) = c) = c$$



Experimental setup:

CIFAR-10

ResNet Wide 32

Accuracy:

Overall: 94%

At 90% confidence: 58%

Accuracy after Temp.Scal:

Overall: 94%

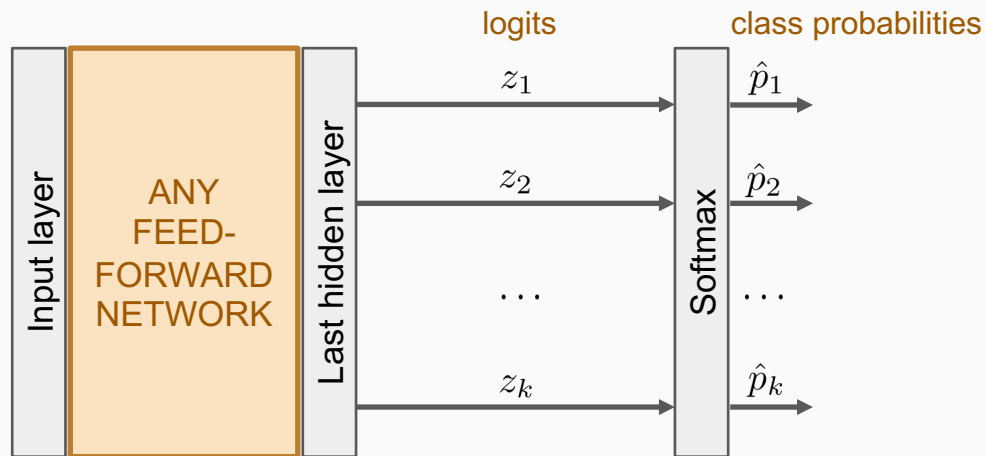
At 90% confidence: 88%

At 90% class 2 prob: 76%

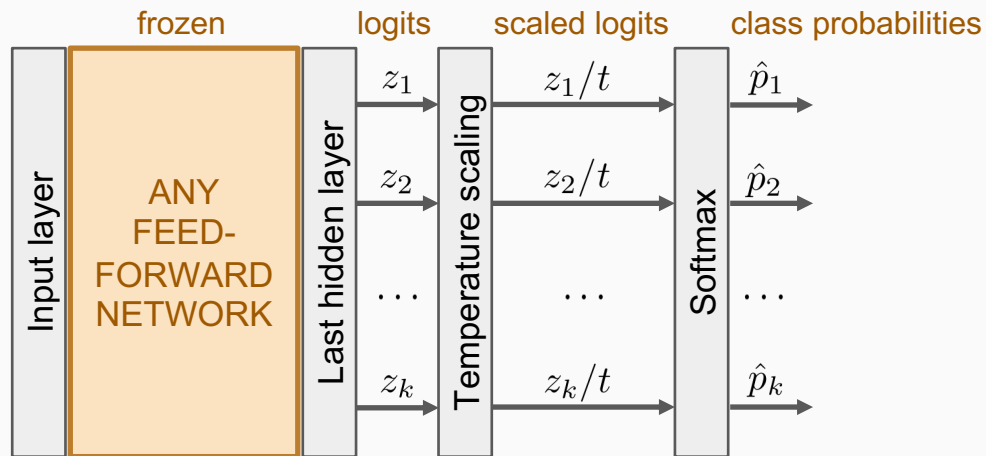
Accuracy after Dir.Calib:

At 90% class 2 prob: 90%

# How to calibrate a multi-class classifier?



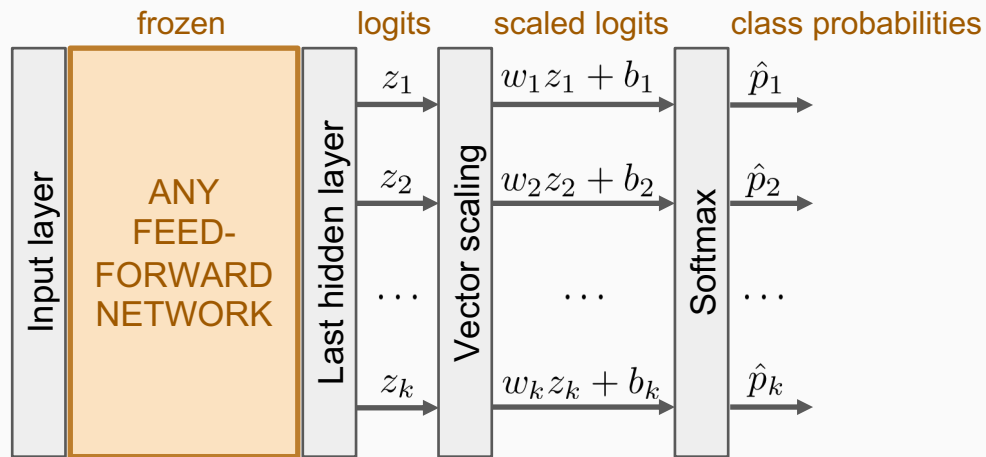
# Temperature scaling



Parameters:  $t \in \mathbb{R}$

C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger. On Calibration of Modern Neural Networks. ICML 2017

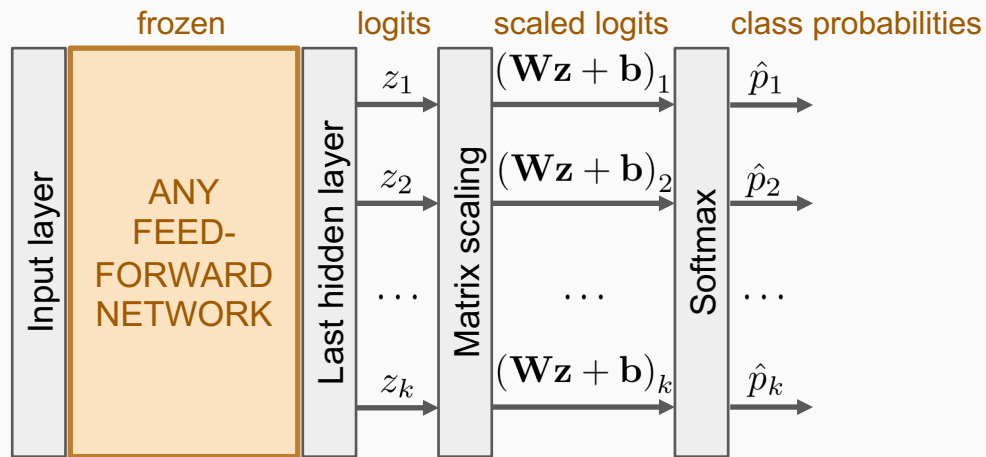
# Vector scaling



Parameters:  $(\mathbf{w}, \mathbf{b}) \in \mathbb{R}^{k+k}$

C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger. On Calibration of Modern Neural Networks. ICML 2017

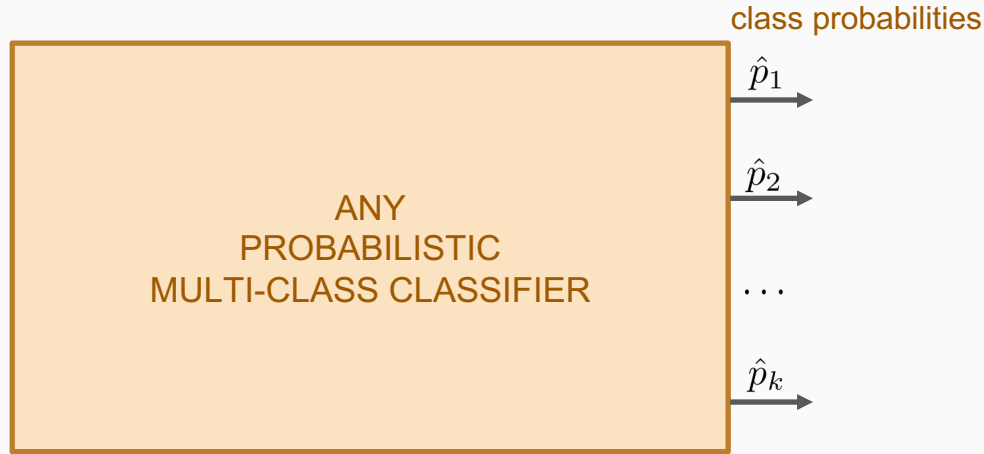
# Matrix scaling



Parameters:  $(\mathbf{W}, \mathbf{b}) \in \mathbb{R}^{k \times k+k}$

C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger. On Calibration of Modern Neural Networks. ICML 2017

# Dirichlet calibration can calibrate any classifiers



# Parametric calibration methods

	Logit space	Class probability space
Binary classification	Derived from Gaussian distribution <b>Platt scaling</b> <sup>[1]</sup>	Derived from Beta distribution <b>Beta calibration</b> <sup>[2]</sup>  (+ constrained variants)
Multi-class classification		

[1] J. Platt. Probabilities for SV machines. In *Advances in Large Margin Classifiers*, pages 61–74, MIT Press, 2000.

[2] M. Kull, T. Silva Filho, P. Flach. Beta calibration: a well-founded and easily implemented improvement on logistic calibration for binary classifiers. AISTATS 2017



UNIVERSITY OF TARTU





# Parametric calibration methods

	Logit space	Class probability space
Binary classification	Derived from Gaussian distribution <b>Platt scaling</b> <sup>[1]</sup>	Derived from Beta distribution <b>Beta calibration</b> <sup>[2]</sup>  (+ constrained variants)
Multi-class classification		Derived from Dirichlet distribution <b>Dirichlet calibration</b>  (+ constrained variants)

[1] J. Platt. Probabilities for SV machines. In *Advances in Large Margin Classifiers*, pages 61–74, MIT Press, 2000.

[2] M. Kull, T. Silva Filho, P. Flach. Beta calibration: a well-founded and easily implemented improvement on logistic calibration for binary classifiers. AISTATS 2017

# Parametric calibration methods

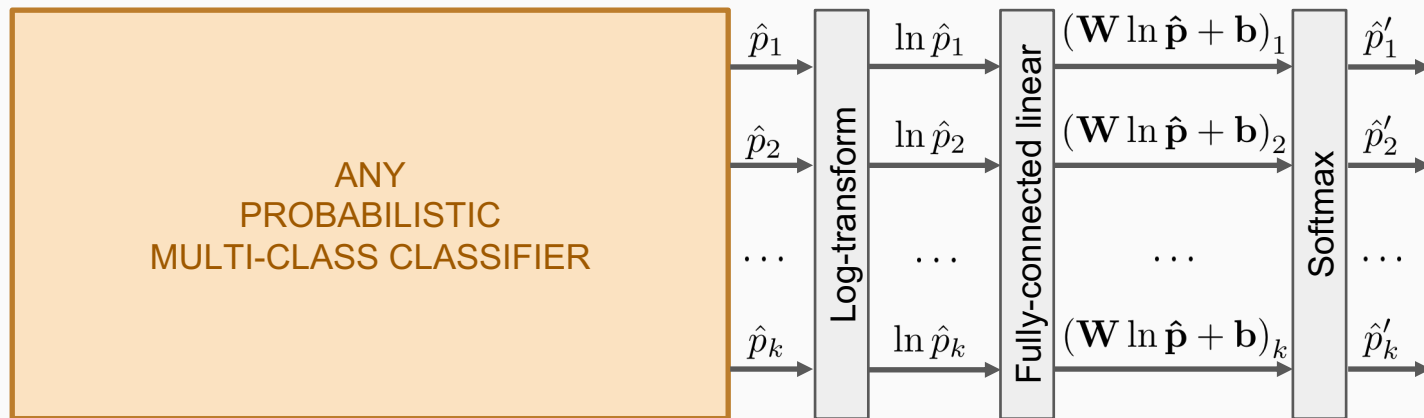
	Logit space	Class probability space
Binary classification	Derived from Gaussian distribution <b>Platt scaling</b> <sup>[1]</sup>	Derived from Beta distribution <b>Beta calibration</b> <sup>[2]</sup>  (+ constrained variants)
Multi-class classification	<b>Matrix scaling</b> <sup>[3]</sup>  (+ vector scaling, temperature scaling)	Derived from Dirichlet distribution <b>Dirichlet calibration</b>  (+ constrained variants)

[1] J. Platt. Probabilities for SV machines. In *Advances in Large Margin Classifiers*, pages 61–74, MIT Press, 2000.

[2] M. Kull, T. Silva Filho, P. Flach. Beta calibration: a well-founded and easily implemented improvement on logistic calibration for binary classifiers. AISTATS 2017

[3] C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger. On Calibration of Modern Neural Networks. ICML 2017

# Dirichlet calibration



Parameters:  $(\mathbf{W}, \mathbf{b}) \in \mathbb{R}^{k \times k+k}$

Regularisation:

- L2
- ODIR (Off-Diagonal and Intercept Regularisation)

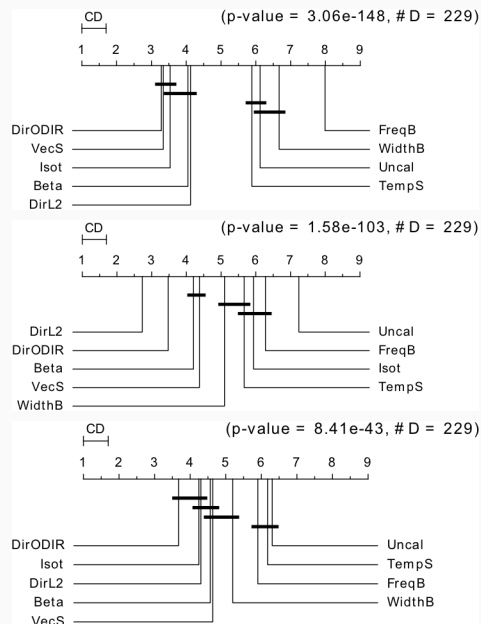
# Non-neural experiments

- 21 datasets x 11 classifiers = 231 settings
- Average rank

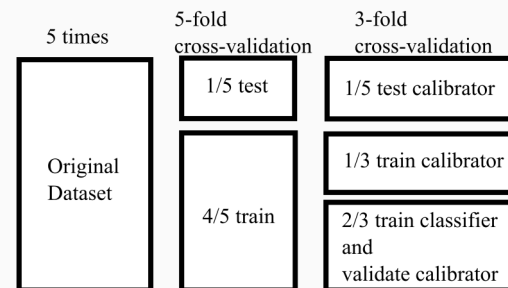
- Classwise-ECE

- Log-loss

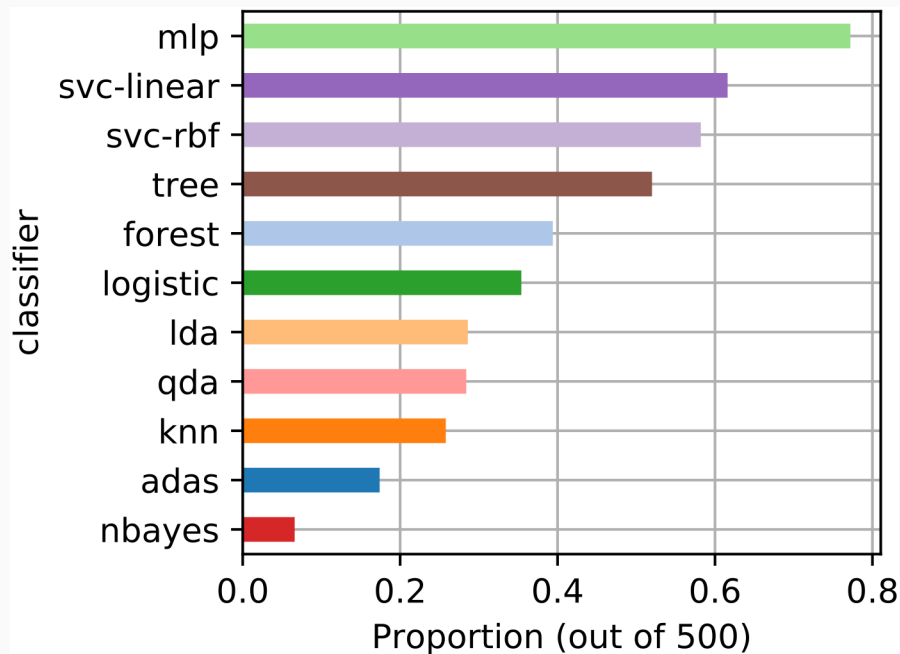
- Error rate



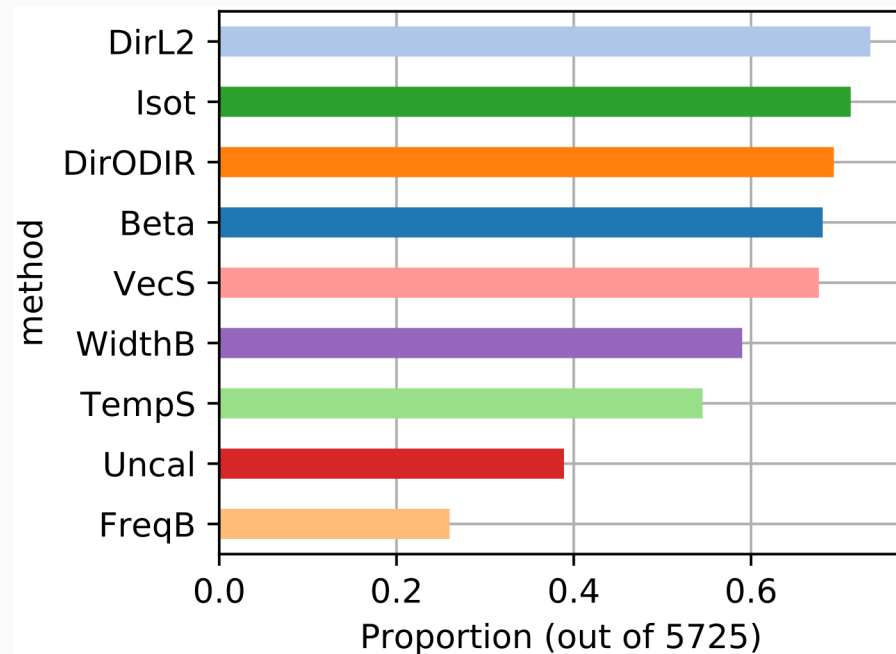
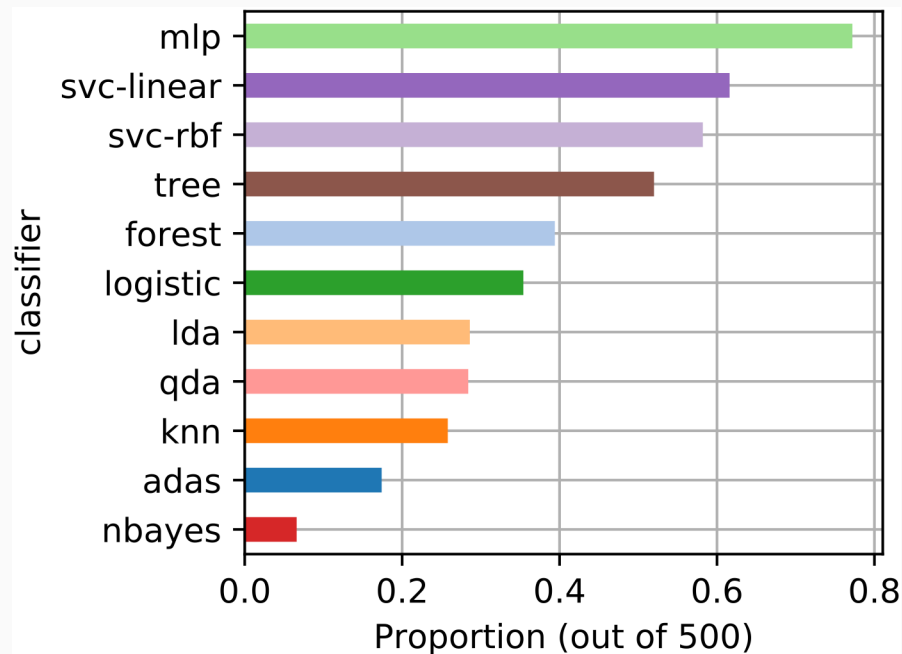
dataset	n_samples	n_features	n_classes
abalone	4177	8	3
balance-scale	625	4	3
car	1728	6	4
cleveland	297	13	5
dermatology	358	34	6
glass	214	9	6
iris	150	4	3
landsat-satellite	6435	36	6
libras-movement	360	90	15
mfeat-karhunen	2000	64	10
mfeat-morphological	2000	6	10
mfeat-zernike	2000	47	10
optdigits	5620	64	10
page-blocks	5473	10	5
pendigits	10992	16	10
segment	2310	19	7
shuttle	101500	9	7
vehicle	846	18	4
vowel	990	10	11
waveform-5000	5000	40	3
yeast	1484	8	10



# Which classifiers are calibrated?



# Which classifiers are calibrated?



# Deep Neural Networks Experiments: Settings

- 3 datasets: CIFAR-10, CIFAR-100, SVHN
- 11 convolutional NNs + 3 pretrained



UNIVERSITY OF TARTU



# Neural experiments

- Datasets: CIFAR-10, CIFAR-100, SVHN
- 11 CNNs trained as in Guo et al + 3 pretrained

Classwise-ECE

	Uncal	general-purpose calibrators			calibrators using logits	
		TempS	Dir-L2	Dir-ODIR	VecS	MS-ODIR
c10_convnet	0.104 <sub>6</sub>	0.044 <sub>4</sub>	0.043 <sub>2</sub>	0.045 <sub>5</sub>	<b>0.043</b> <sub>1</sub>	0.044 <sub>3</sub>
c10_densenet40	0.114 <sub>6</sub>	0.040 <sub>5</sub>	<b>0.034</b> <sub>1</sub>	0.037 <sub>4</sub>	0.036 <sub>2</sub>	0.037 <sub>3</sub>
c10_lenet5	0.198 <sub>6</sub>	0.171 <sub>5</sub>	<b>0.052</b> <sub>1</sub>	0.059 <sub>4</sub>	0.057 <sub>2</sub>	0.059 <sub>3</sub>
c10_resnet110	0.098 <sub>6</sub>	0.043 <sub>5</sub>	<b>0.032</b> <sub>1</sub>	0.039 <sub>4</sub>	0.037 <sub>3</sub>	0.036 <sub>2</sub>
c10_resnet110_SD	0.086 <sub>6</sub>	0.031 <sub>4</sub>	0.031 <sub>5</sub>	0.029 <sub>3</sub>	0.027 <sub>2</sub>	<b>0.027</b> <sub>1</sub>
c10_resnet_wide32	0.095 <sub>6</sub>	0.048 <sub>5</sub>	0.032 <sub>3</sub>	0.029 <sub>2</sub>	0.032 <sub>4</sub>	<b>0.029</b> <sub>1</sub>
c100_convnet	0.424 <sub>6</sub>	<b>0.227</b> <sub>1</sub>	0.402 <sub>5</sub>	0.240 <sub>3</sub>	0.241 <sub>4</sub>	0.240 <sub>2</sub>
c100_densenet40	0.470 <sub>6</sub>	0.187 <sub>2</sub>	0.330 <sub>5</sub>	<b>0.186</b> <sub>1</sub>	0.189 <sub>3</sub>	0.191 <sub>4</sub>
c100_lenet5	0.473 <sub>6</sub>	0.385 <sub>5</sub>	0.219 <sub>4</sub>	0.213 <sub>2</sub>	<b>0.203</b> <sub>1</sub>	0.214 <sub>3</sub>
c100_resnet110	0.416 <sub>6</sub>	0.201 <sub>3</sub>	0.359 <sub>5</sub>	<b>0.186</b> <sub>1</sub>	0.194 <sub>2</sub>	0.203 <sub>4</sub>
c100_resnet110_SD	0.375 <sub>6</sub>	0.203 <sub>4</sub>	0.373 <sub>5</sub>	0.189 <sub>3</sub>	<b>0.170</b> <sub>1</sub>	0.186 <sub>2</sub>
c100_resnet_wide32	0.420 <sub>6</sub>	0.186 <sub>4</sub>	0.333 <sub>5</sub>	0.180 <sub>2</sub>	<b>0.171</b> <sub>1</sub>	0.180 <sub>3</sub>
SVHN_convnet	0.159 <sub>6</sub>	0.038 <sub>4</sub>	0.043 <sub>5</sub>	0.026 <sub>2</sub>	<b>0.025</b> <sub>1</sub>	0.027 <sub>3</sub>
SVHN_resnet152_SD	0.019 <sub>2</sub>	<b>0.018</b> <sub>1</sub>	0.022 <sub>6</sub>	0.020 <sub>3</sub>	0.021 <sub>5</sub>	0.021 <sub>4</sub>
Average rank	5.71	3.71	3.79	2.79	2.29	2.71

Log-loss

	Uncal	general-purpose calibrators			calibrators using logits	
		TempS	Dir-L2	Dir-ODIR	VecS	MS-ODIR
c10_convnet	0.391 <sub>6</sub>	<b>0.195</b> <sub>1</sub>	0.197 <sub>4</sub>	0.195 <sub>2</sub>	0.197 <sub>5</sub>	0.196 <sub>3</sub>
c10_densenet40	0.428 <sub>6</sub>	0.225 <sub>5</sub>	<b>0.220</b> <sub>1</sub>	0.224 <sub>4</sub>	0.223 <sub>3</sub>	0.222 <sub>2</sub>
c10_lenet5	0.823 <sub>6</sub>	0.800 <sub>5</sub>	0.744 <sub>2</sub>	0.744 <sub>3</sub>	0.747 <sub>4</sub>	<b>0.743</b> <sub>1</sub>
c10_resnet110	0.358 <sub>6</sub>	0.209 <sub>5</sub>	<b>0.203</b> <sub>1</sub>	0.205 <sub>3</sub>	0.206 <sub>4</sub>	0.204 <sub>2</sub>
c10_resnet110_SD	0.303 <sub>6</sub>	0.178 <sub>5</sub>	0.177 <sub>4</sub>	0.176 <sub>3</sub>	0.175 <sub>2</sub>	<b>0.175</b> <sub>1</sub>
c10_resnet_wide32	0.382 <sub>6</sub>	0.191 <sub>5</sub>	0.185 <sub>4</sub>	0.182 <sub>2</sub>	0.183 <sub>3</sub>	<b>0.182</b> <sub>1</sub>
c100_convnet	1.641 <sub>6</sub>	<b>0.942</b> <sub>1</sub>	1.189 <sub>5</sub>	0.961 <sub>2</sub>	0.964 <sub>4</sub>	0.961 <sub>3</sub>
c100_densenet40	2.017 <sub>6</sub>	1.057 <sub>2</sub>	1.253 <sub>5</sub>	1.059 <sub>4</sub>	1.058 <sub>3</sub>	<b>1.051</b> <sub>1</sub>
c100_lenet5	2.784 <sub>6</sub>	2.650 <sub>5</sub>	2.595 <sub>4</sub>	2.490 <sub>2</sub>	2.516 <sub>3</sub>	<b>2.487</b> <sub>1</sub>
c100_resnet110	1.694 <sub>6</sub>	1.092 <sub>3</sub>	1.212 <sub>5</sub>	1.096 <sub>4</sub>	1.089 <sub>2</sub>	<b>1.074</b> <sub>1</sub>
c100_resnet110_SD	1.353 <sub>6</sub>	0.942 <sub>3</sub>	1.198 <sub>5</sub>	0.945 <sub>4</sub>	<b>0.923</b> <sub>1</sub>	0.927 <sub>2</sub>
c100_resnet_wide32	1.802 <sub>6</sub>	0.945 <sub>3</sub>	1.087 <sub>5</sub>	0.953 <sub>4</sub>	0.937 <sub>2</sub>	<b>0.933</b> <sub>1</sub>
SVHN_convnet	0.205 <sub>6</sub>	0.151 <sub>5</sub>	0.142 <sub>3</sub>	0.138 <sub>2</sub>	0.144 <sub>4</sub>	<b>0.138</b> <sub>1</sub>
SVHN_resnet152_SD	0.085 <sub>6</sub>	<b>0.079</b> <sub>1</sub>	0.085 <sub>5</sub>	0.080 <sub>2</sub>	0.081 <sub>4</sub>	0.081 <sub>3</sub>
Average rank	6.0	3.5	3.79	2.93	3.14	1.64



# Conclusion

1. Dirichlet calibration: New parametric general-purpose multiclass calibration method
  - a. Natural extension of two-class Beta calibration
  - b. Easy to implement with multinomial logistic regression on log-transformed class probabilities
2. Best or tied best performance with 21 datasets x 11 classifiers
3. Advances state-of-the-art on Neural Networks by introducing ODIR regularisation



UNIVERSITY OF TARTU



# Beyond temperature scaling: Obtaining well-calibrated multiclass probabilities with Dirichlet calibration

Meelis Kull, Miquel Perello Nieto, Markus Kängsepp,  
Telmo Silva Filho, Hao Song, Peter Flach

NeurIPS 2019



UNIVERSITY OF TARTU



University of  
BRISTOL



Universidade Federal da Paraíba  
ESTATÍSTICA



The  
Alan Turing  
Institute

