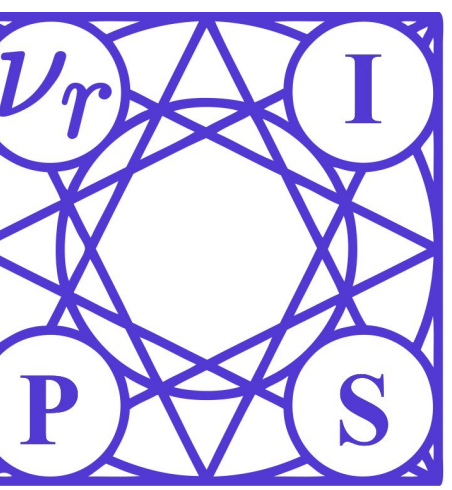


# Beyond temperature scaling:

## Obtaining well-calibrated multiclass probabilities with Dirichlet calibration

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



Class probabilities predicted by most multiclass classifiers are uncalibrated, often tending towards over-confidence. With neural networks, calibration can be improved by temperature scaling, a method to learn a single corrective multiplicative factor for inputs to the last softmax layer. On non-neural models the existing methods apply binary calibration in a pairwise or one-vs-rest fashion. We propose a natively multiclass calibration method applicable to classifiers from any model class, derived from Dirichlet distributions and generalising the beta calibration method from binary classification. It is easily implemented with neural nets since it is equivalent to log-transforming the uncalibrated probabilities, followed by one linear layer and softmax.

### Notion of calibration:

- A probabilistic classifier  $\hat{\mathbf{p}}$  is multiclass-calibrated if for any prediction vector  $\hat{\mathbf{p}}(\mathbf{x}) = \mathbf{q}$ , the proportions of classes among all possible instances  $\mathbf{x}$  getting the same prediction are:  $P(Y = i \mid \hat{\mathbf{p}}_i(X) = \mathbf{q}) = q_i$  for  $i = 1, \dots, k$ .

- Classwise-calibrated if:  $P(Y = i \mid \hat{\mathbf{p}}_i(X) = q_i) = q_i$ .

- Confidence-calibrated if:  $P(Y = \arg \max(\hat{\mathbf{p}}(X)) \mid \max(\hat{\mathbf{p}}(X)) = c) = c$ .

### Evaluation of calibration:

- Some empirical measures of calibration include:

$$\text{confidence-ECE} = \frac{1}{n} \sum_{i=1}^n |y_j(B_i) - \hat{p}_j(B_i)|$$

- Other measures include well-known proper scoring rules, such as Brier score and log-loss

- Every proper loss is minimised by the canonical calibration function

### Temperature scaling:

- Temperature scaling learns a temperature parameter  $t > 0$  which decreases ( $t > 1$ ) or increases ( $t < 1$ ) the confidence of a model with softmax output

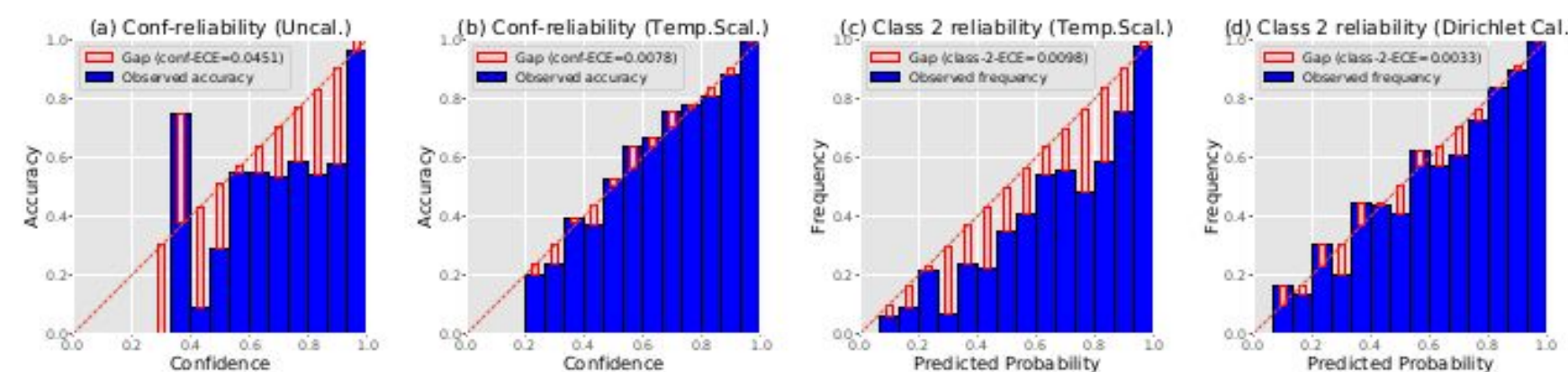


Figure 1: Reliability diagrams of c10\_resnet\_wide32 on CIFAR-10: (a) confidence-reliability before calibration; (b) confidence-reliability after temperature scaling; (c) classwise-reliability for class 2 after temperature scaling; (d) classwise-reliability for class 2 after Dirichlet calibration.

### Dirichlet calibration:

- We consider the distribution of prediction vectors  $\hat{\mathbf{p}}(\mathbf{x})$  separately on instances of each class, and assume these are Dirichlet distributions with different parameters:

$$\hat{\mathbf{p}}(X) \mid Y = j \sim \text{Dir}(\alpha^{(j)})$$

$$\text{generative parametrisation: } \hat{\mu}_{\text{DirGen}}(\mathbf{q}; \alpha, \pi) = (\pi_1 f_1(\mathbf{q}), \dots, \pi_k f_k(\mathbf{q})) / z$$

$$\text{linear parametrisation: } \hat{\mu}_{\text{DirLin}}(\mathbf{q}; \mathbf{W}, \mathbf{b}) = \sigma(\mathbf{W} \ln \mathbf{q} + \mathbf{b})$$

$$\text{canonical parametrisation: } \hat{\mu}_{\text{Dir}}(\mathbf{q}; \mathbf{A}, \mathbf{c}) = \sigma(\mathbf{A} \ln \frac{\mathbf{q}}{1/k} + \ln \mathbf{c})$$

- All parametrisations are equal, i.e. they contain exactly the same calibration maps.

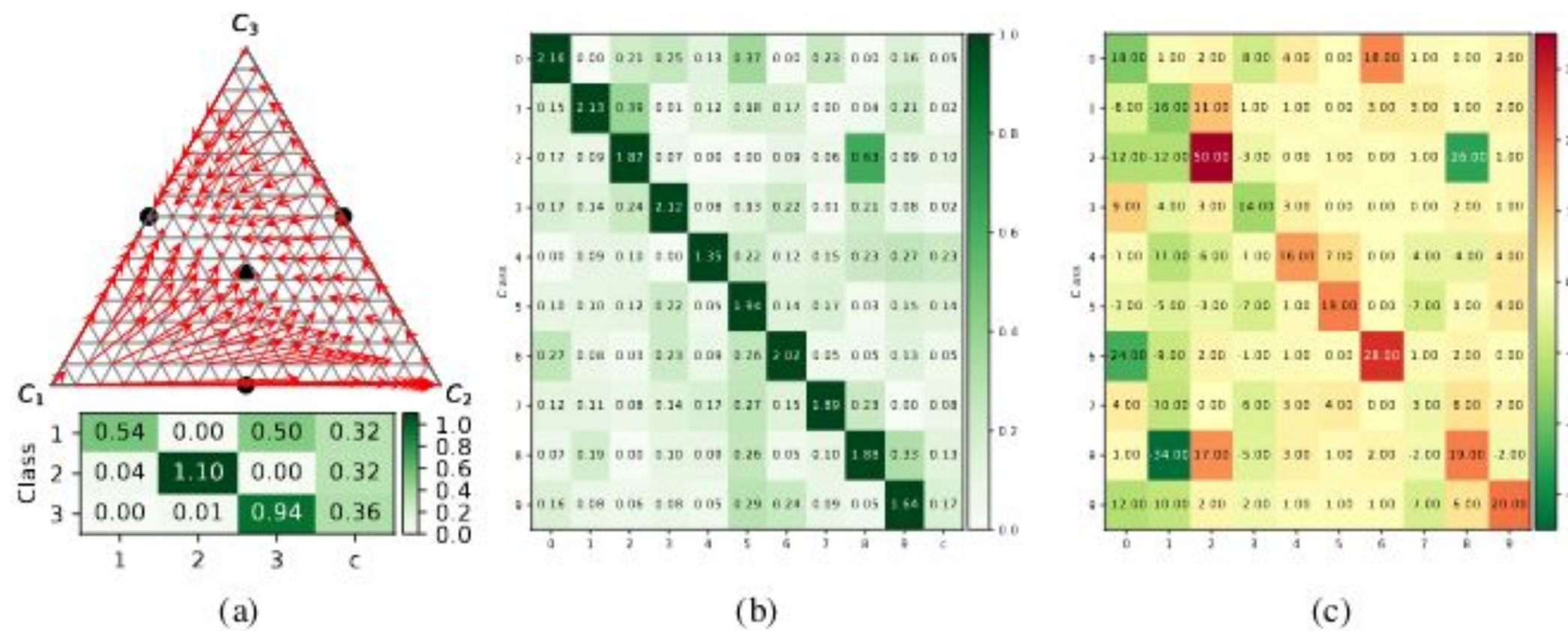
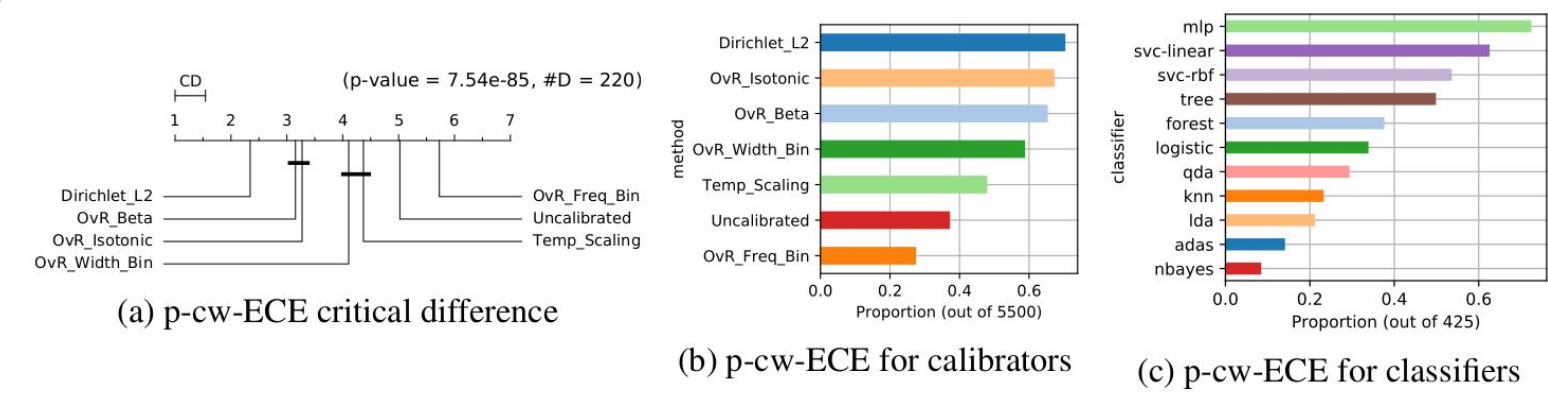


Figure 2: Interpretation of Dirichlet calibration maps: (a) calibration map for MLP on the abalone dataset, 4 interpretation points shown by black dots, and canonical parametrisation as a matrix with  $\mathbf{A}, \mathbf{c}$ ; (b) canonical parametrisation of a map on SVHN\_convnet; (c) changes to the confusion matrix after applying this calibration map.

### Non-neural experiment:

- 21 datasets and 11 classifiers = 231 settings
- 6 calibration methods:
- Best or tied best performance



### Deep Neural experiment:

- 3 datasets CIFAR-10, CIFAR-100, SVHN
- 11 convolutional NNs
- 3 pretrained CNNs
- 5 calibration methods:
- 8 evaluation measures
- 5-fold-crossval.

Table 3: Scores and ranking of calibration methods for cw-ECE.

	Uncal	general-purpose calibrators		calibrators using logits			
		Temps	Dir-L2	Dir-ODIR	Yes	No	MS-ODIR
c10_convnet	0.104 <sub>1</sub>	0.044 <sub>1</sub>	0.043 <sub>1</sub>	0.045 <sub>1</sub>	<b>0.043<sub>1</sub></b>	0.044 <sub>1</sub>	
c10_densenet40	0.114 <sub>1</sub>	0.040 <sub>1</sub>	<b>0.034<sub>1</sub></b>	0.037 <sub>1</sub>	0.036 <sub>2</sub>	0.037 <sub>1</sub>	
c10_inet5	0.198 <sub>1</sub>	0.171 <sub>1</sub>	<b>0.062<sub>1</sub></b>	0.059 <sub>1</sub>	0.057 <sub>2</sub>	0.059 <sub>1</sub>	
c10_resnet10	0.098 <sub>1</sub>	0.043 <sub>1</sub>	<b>0.032<sub>1</sub></b>	0.039 <sub>1</sub>	0.037 <sub>2</sub>	0.036 <sub>2</sub>	
c10_resnet10_SD	0.088 <sub>1</sub>	0.031 <sub>1</sub>	0.031 <sub>1</sub>	0.029 <sub>1</sub>	0.027 <sub>2</sub>	<b>0.027<sub>1</sub></b>	
c10_resnet_wide32	0.055 <sub>1</sub>	0.045 <sub>1</sub>	0.052 <sub>1</sub>	0.029 <sub>1</sub>	0.032 <sub>2</sub>	<b>0.027<sub>1</sub></b>	
c100_convnet	0.424 <sub>1</sub>	<b>0.227<sub>1</sub></b>	0.402 <sub>1</sub>	0.240 <sub>1</sub>	0.241 <sub>2</sub>	0.240 <sub>2</sub>	
c100_densenet40	0.474 <sub>1</sub>	0.187 <sub>1</sub>	0.330 <sub>1</sub>	<b>0.186<sub>1</sub></b>	0.189 <sub>1</sub>	0.191 <sub>1</sub>	
c100_inet5	0.473 <sub>1</sub>	0.385 <sub>1</sub>	0.219 <sub>1</sub>	0.213 <sub>1</sub>	<b>0.203<sub>1</sub></b>	0.214 <sub>1</sub>	
c100_resnet10	0.416 <sub>1</sub>	0.201 <sub>1</sub>	0.359 <sub>1</sub>	<b>0.186<sub>1</sub></b>	0.194 <sub>1</sub>	0.203 <sub>1</sub>	
c100_resnet10_SD	0.375 <sub>1</sub>	0.203 <sub>1</sub>	0.375 <sub>1</sub>	0.189 <sub>1</sub>	<b>0.179<sub>1</sub></b>	0.186 <sub>2</sub>	
c100_resnet_wide32	0.428 <sub>1</sub>	0.186 <sub>1</sub>	0.352 <sub>1</sub>	<b>0.189<sub>1</sub></b>	<b>0.173<sub>1</sub></b>	0.186 <sub>2</sub>	
SVHN_convnet	0.159 <sub>1</sub>	0.038 <sub>1</sub>	0.043 <sub>1</sub>	0.036 <sub>1</sub>	<b>0.025<sub>1</sub></b>	0.027 <sub>1</sub>	
SVHN_resnet52_SD	0.019 <sub>1</sub>	<b>0.018<sub>1</sub></b>	0.022 <sub>1</sub>	0.020 <sub>1</sub>	0.021 <sub>2</sub>	0.021 <sub>1</sub>	
Average rank	5.71	3.71	3.79	2.79	2.29	2.71	

Table 4: Scores and ranking of calibration methods for log-loss.

	Uncal	general-purpose calibrators			calibrators using logits	
		Temps	Dir-L2	Dir-ODIR	VecS	MS-ODIR
c10_convnet	0.391 <sub>1</sub>	0.195 <sub>1</sub>	0.197 <sub>1</sub>	0.195 <sub>1</sub>	0.197 <sub>1</sub>	0.196 <sub>1</sub>
c10_densenet40	0.428 <sub>1</sub>	0.225 <sub>1</sub>	<b>0.220<sub>1</sub></b>	0.224 <sub>1</sub>	0.223 <sub>1</sub>	0.222 <sub>1</sub>
c10_inet5	0.825 <sub>1</sub>	0.808 <sub>1</sub>	0.744 <sub>1</sub>	0.744 <sub>1</sub>	0.747 <sub>1</sub>	0.743 <sub>1</sub>
c10_resnet10	0.358 <sub>1</sub>	0.206 <sub>1</sub>	<b>0.203<sub>1</sub></b>	0.205 <sub>1</sub>	0.206 <sub>1</sub>	0.204 <sub>1</sub>
c10_resnet10_SD	0.303 <sub>1</sub>	0.178 <sub>1</sub>	0.177 <sub>1</sub>	0.176 <sub>1</sub>	0.175 <sub>1</sub>	<b>0.175<sub>1</sub></b>
c10_resnet_wide32	0.382 <sub>1</sub>	0.191 <sub>1</sub>	0.185 <sub>1</sub>	0.182 <sub>1</sub>	0.183 <sub>1</sub>	<b>0.182<sub>1</sub></b>
c100_convnet	1.641 <sub>1</sub>	<b>0.942<sub>1</sub></b>	1.189 <sub>1</sub>	0.961 <sub>1</sub>	0.964 <sub>1</sub>	0.961 <sub>1</sub>
c100_densenet40	2.017 <sub>1</sub>	1.075 <sub>1</sub>	1.255 <sub>1</sub>	1.059 <sub>1</sub>	1.058 <sub>1</sub>	<b>1.051<sub>1</sub></b>
c100_inet5	2.784 <sub>1</sub>	2.650 <sub>1</sub>	2.595 <sub>1</sub>	2.490 <sub>1</sub>	2.516 <sub>1</sub>	2.487 <sub>1</sub>
c100_resnet10	1.664 <sub>1</sub>	1.092 <sub>1</sub>	1.212 <sub>1</sub>	1.096 <sub>1</sub>	1.089 <sub>1</sub>	<b>1.074<sub>1</sub></b>
c100_resnet10_SD	1.353 <sub>1</sub>	0.942 <sub>1</sub>	1.186 <sub>1</sub>	0.945 <sub>1</sub>	<b>0.921<sub>1</sub></b>	0.927 <sub>1</sub>
c100_resnet_wide32	1.862 <sub>1</sub>	0.945 <sub>1</sub>	1.087 <sub>1</sub>	0.952 <sub>1</sub>	0.937 <sub>1</sub>	<b>0.933<sub>1</sub></b>
SVHN_convnet	0.206 <sub>1</sub>	0.151 <sub>1</sub>	0.142 <sub>1</sub>	0.138 <sub>1</sub>	0.144 <sub>1</sub>	<b>0.138<sub>1</sub></b>
SVHN_resnet52_SD	0.085 <sub>1</sub>	<b>0.079<sub>1</sub></b>	0.085 <sub>1</sub>	0.080 <sub>1</sub>	0.081 <sub>1</sub>	0.081 <sub>1</sub>
Average rank	6.0	3.5	3.79	2.93	3.14	1.64

### Conclusion:

- Dirichlet calibration: New parametric general-purpose multiclass calibration method
  - Natural extension of two-class Beta calibration
  - Easy to implement with multinomial logistic regression on log-transformed class probabilities
- Advances state-of-the-art on Neural Networks by introducing ODIR regularisation

### Future work:

- Which neural architectures and train have temperature scaling as calibration function
- Use other distributions of the exponential family
- Investigate scores coming from distributions per class

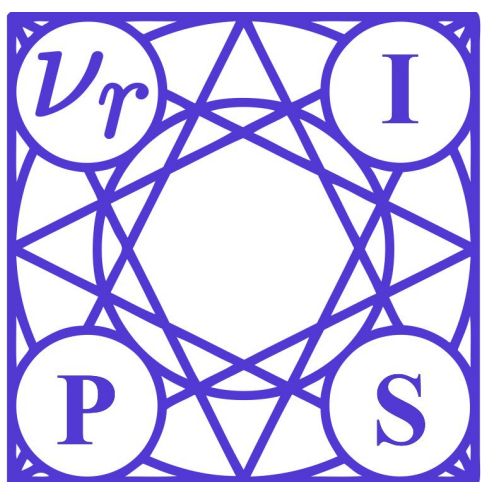




Beyond temperature scaling:

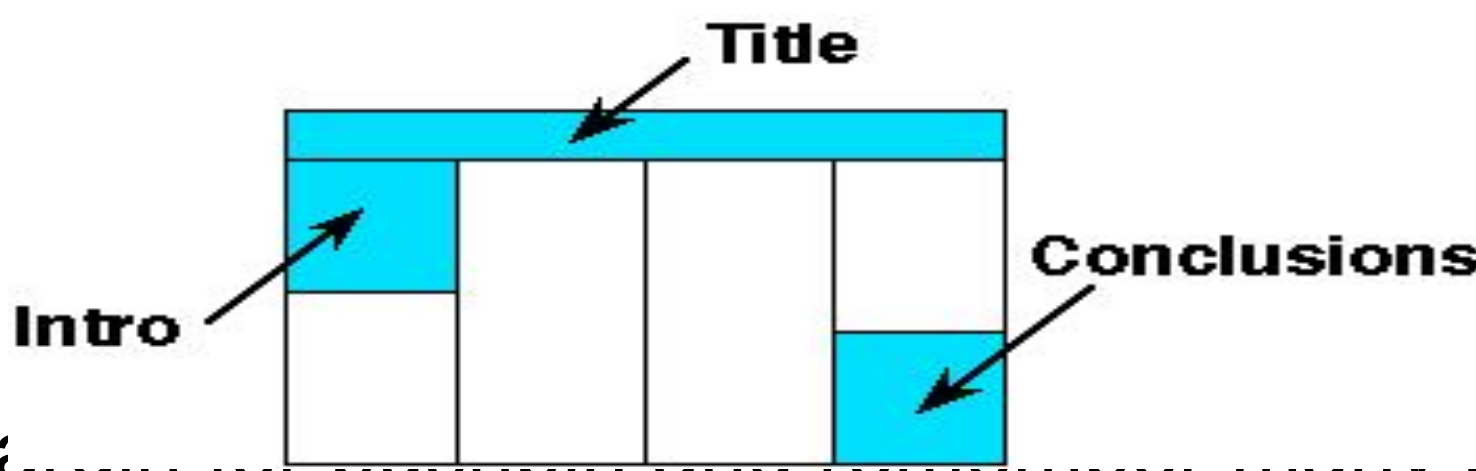
Obtaining well-calibrated multiclass probabilities with Dirichlet calibration

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



### Section 1 (sizes):

- Posters boards are 48” tall and 96” wide, but we recommend you leave a little border since you may not be able to pin at the vertical edge. Since PowerPoint does not let one define such a large paper size, this template is designed to be printed at 200%, yielding a 46” x94” poster. You can scale it up or down a bit (e.g. 42” is a common paper size at FexEd). Note there is no direct international A0.. A1 equivalent. The poster size is approximately three A0 boards next to each other, i.e., each column in this example is about one A0 board.
- Ideally you want to keep it very readable: this is not your paper, it is a poster. 32pt here (64 final printing) is good for most text:
  - Sub-bullets are 28 here (56 final)
  - Don’t use smaller than 24pt in this template (which is 48pt in final printing at 200%)
  - Insert plenty of graphics and any math you need
- When inserting graphics or equations, keep the resolution high (remember this will be printed at 200%). If you can see blocking artifacts at 400% magnification in PowerPoint, consider finding better graphics. This is an example of BAD/LOW RES GRAPHICS



- Leave enough margin for captions and remember many big plotters cannot get within .5” of the actual paper edge.
- You are free to use colored backgrounds and such but they generally reduce readability.
- You are free to use what ever fonts you like.
  - San Serif fonts like Arial are more readable from a distance,
  - Serif fonts like times may look more consistent with your mathematics

### Section 2 (layout):

- Remember the poster session will be crowded so design the poster to be read in columns so people can read what is in front of them and move left to right to get the whole story.
- The poster should use photos, figures, and tables to tell the story of the study. For clarity, present the information in a sequence that is easy to follow.
- There is often way too much text in a poster - there definitely is in this template! Posters primarily are visual presentations; the text should support the graphics. Look critically at the layout. Some poster 'experts' suggest that if there is about 20-25% text, 40-45% graphics and 30-40% empty space, you are doing well.

### Section 3:

- Include more figures than are in the paper so you can talk to them. Include things that are not in the paper and then encourage them to read the paper. Don’t try to just put all the paper here.
- If it looks like a cut/paste of the paper, people skip that poster since they can read the papers after the conference. Many people find it better to spend time talking with poster presenters that have more to offer than just redoing the paper content paper in big fonts.
- Remember Poster boards look like this.. This is your canvas. Paint us a picture of your work.



Maybe a QrCode to the website with your code.



