

A. Proofs

Proof of Proposition 3.1 The proof is analogous to [30]. The L_1 calibration error for binary classification task is given by $CE_{cls} := \mathbb{E}[|\mathbb{E}[Y = 1 | \hat{S}] - \hat{S}|]$.

Proof.

$$d_{cls} = |\mathbb{E}[Y - \hat{S}]| \quad (16)$$

$$= |\mathbb{E}[Y - \mathbb{E}[Y = 1 | \hat{S}]] + \mathbb{E}[\mathbb{E}[Y = 1 | \hat{S}] - \hat{S}]| \quad (17)$$

$$= \mathbb{E}[Y] - \underbrace{\mathbb{E}[\mathbb{E}[Y = 1 | \hat{S}]]}_{=\mathbb{E}[Y]} + \mathbb{E}[\mathbb{E}[Y = 1 | \hat{S}] - \hat{S}] \quad (18)$$

$$\leq \underbrace{\mathbb{E}[|\mathbb{E}[Y = 1 | \hat{S}] - \hat{S}|]}_{=CE_{cls}}. \quad (19)$$

The inequality in (19) is obtained by utilizing the convexity of the absolute value and applying Jensen’s inequality. \square

Proof of Proposition 3.2

Proof.

$$CE_{det} = \mathbb{E}[|\mathbb{E}[IoU | \hat{S}] - \hat{S}|] \quad (20)$$

$$= \mathbb{E}[|\mathbb{E}[IoU | \hat{S}] - \mathbb{E}[\hat{S} | \hat{S}]|] \quad (21)$$

$$\leq \mathbb{E}[|\mathbb{E}[IoU - \hat{S} | \hat{S}]|] \quad (22)$$

$$= \mathbb{E}[|IoU - \hat{S}|] = d_{det} \quad (23)$$

\square

As before, (22) is obtained by applying Jensen’s inequality.

B. Experiments

In this section we will provide additional quantitative and qualitative results, following the structure of the experiments in the main text.

We compare the performance of the binning-based LaECE metric with our estimator \widehat{CE} , for measuring calibration error as obtained by using an identity link in Equation (4). In Figure 4 we can observe that both metrics perform equally well and converge to the ground truth value after 5000 points.

In Tables 7 and 8, we investigate the performance of our estimator in a calibration regularized training framework on Cityscapes and Pascal VOC datasets, and compare it with competing post-hoc and trainable methods. We report two versions of AP, \widehat{CE} , and D-ECE: averaged over 10 thresholds and evaluated at IoU = 0.5. Additionally, we include \widehat{CE} with an identity link. The models with \widehat{CE} and TCD

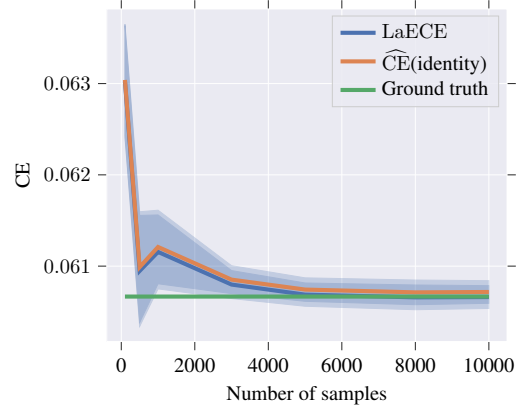


Figure 4. Comparison of LaECE vs. \widehat{CE} (identity) as a function of the number of points used for the estimation.

auxiliary loss are trained from scratch. We show the results for several values of λ and observe that by increasing the weight of the regularization, the calibration error is reduced, as expected. In all evaluations we use a detection score threshold $\gamma = 0.5$, and the results are averaged across three seeds.

In Tables 9 and 10 we summarize the AP and CE metrics for fine-tuning on Cityscapes and COCO, respectively. The learning rate for fine-tuning is set to the one from the last training iteration, i.e., $1e-3$ for Cityscapes and $1e-4$ for COCO and Pascal VOC.

Figure 5 shows qualitative results on COCO, comparing uncalibrated F-RCNN model with a fine-tuned model using \widehat{CE} loss. Our estimator effectively adjusts the confidence scores to achieve better alignment with the IoU overlap, thus also reducing the number of wrong detections.

Table 7. Calibration and detection performance of models trained on Cityscapes.

Model	AP	AP ₅₀	\widehat{CE}	\widehat{CE}_{50}	\widehat{CE} (identity)	D-ECE	D-ECE ₅₀
F-RCNN	35.36 \pm 0.65	54.46 \pm 0.82	37.45 \pm 0.46	17.00 \pm 0.46	28.48 \pm 0.68	37.83 \pm 0.76	16.65 \pm 0.91
F-RCNN + TS	32.29 \pm 0.25	49.55 \pm 0.27	26.36 \pm 0.86	13.59 \pm 0.33	11.90 \pm 1.22	26.05 \pm 0.62	13.71 \pm 0.48
F-RCNN + \widehat{CE} ($\lambda = 1$)	35.12 \pm 0.42	54.52 \pm 0.46	36.62 \pm 0.90	16.06 \pm 1.22	27.08 \pm 1.17	36.41 \pm 0.66	14.96 \pm 0.96
F-RCNN + \widehat{CE} ($\lambda = 2$)	35.46 \pm 0.07	53.58 \pm 0.32	36.61 \pm 0.45	16.24 \pm 0.32	27.62 \pm 0.74	37.20 \pm 0.44	16.33 \pm 0.22
F-RCNN + \widehat{CE} ($\lambda = 3$)	35.04 \pm 0.26	53.91 \pm 0.12	36.99 \pm 0.33	16.61 \pm 0.27	27.30 \pm 0.69	36.84 \pm 0.60	15.64 \pm 0.39
F-RCNN + \widehat{CE} ($\lambda = 4$)	35.45 \pm 0.30	53.97 \pm 0.05	35.93 \pm 0.75	15.33 \pm 0.79	26.43 \pm 1.23	36.34 \pm 0.89	15.24 \pm 1.01
RetinaNet	34.60 \pm 0.23	52.60 \pm 0.81	32.22 \pm 0.29	12.91 \pm 0.95	23.25 \pm 0.39	32.51 \pm 0.18	12.54 \pm 0.66
RetinaNet + TS	34.60 \pm 0.23	52.60 \pm 0.81	28.72 \pm 1.25	11.63 \pm 1.19	18.67 \pm 2.07	28.92 \pm 1.18	11.05 \pm 1.25
RetinaNet + TCD	33.94 \pm 0.81	51.83 \pm 1.44	33.79 \pm 0.99	13.35 \pm 0.88	25.01 \pm 1.13	33.90 \pm 0.87	13.04 \pm 0.47
RetinaNet + \widehat{CE} ($\lambda = 0.05$)	32.59 \pm 0.71	49.22 \pm 1.38	28.77 \pm 1.94	11.90 \pm 0.19	17.33 \pm 3.15	28.55 \pm 2.09	11.07 \pm 0.22
RetinaNet + \widehat{CE} ($\lambda = 0.1$)	33.03 \pm 0.39	50.95 \pm 0.75	30.21 \pm 0.45	10.90 \pm 0.28	20.15 \pm 0.11	30.18 \pm 0.49	10.30 \pm 0.35
FCOS	34.81 \pm 0.08	52.31 \pm 0.29	26.29 \pm 0.98	13.91 \pm 1.28	14.43 \pm 1.31	25.83 \pm 0.98	13.23 \pm 1.18
FCOS + TS	33.62 \pm 0.44	50.23 \pm 1.17	26.45 \pm 1.23	13.06 \pm 0.86	13.33 \pm 2.16	25.46 \pm 1.38	11.96 \pm 0.90
FCOS + TCD	35.57 \pm 0.23	54.21 \pm 0.38	28.53 \pm 0.82	15.09 \pm 0.83	16.85 \pm 1.27	27.78 \pm 0.77	13.68 \pm 0.86
FCOS + \widehat{CE} ($\lambda = 0.05$)	34.74 \pm 0.55	52.80 \pm 0.69	26.39 \pm 0.48	15.27 \pm 1.08	15.20 \pm 0.75	26.35 \pm 0.54	14.78 \pm 0.98
FCOS + \widehat{CE} ($\lambda = 0.1$)	34.47 \pm 0.52	52.34 \pm 0.71	26.88 \pm 0.33	16.73 \pm 1.46	15.63 \pm 0.47	26.74 \pm 0.09	15.75 \pm 1.71
FCOS + \widehat{CE} ($\lambda = 0.5$)	33.73 \pm 0.57	51.23 \pm 1.12	25.24 \pm 0.85	18.15 \pm 1.07	13.07 \pm 0.83	24.60 \pm 0.49	17.45 \pm 0.61
FCOS + \widehat{CE} ($\lambda = 1$)	32.91 \pm 0.61	50.09 \pm 0.66	26.30 \pm 0.28	19.88 \pm 0.92	14.73 \pm 0.63	25.80 \pm 0.31	18.80 \pm 0.79

Table 8. Calibration and detection performance of models trained on Pascal VOC.

Model	AP	AP ₅₀	\widehat{CE}	\widehat{CE}_{50}	\widehat{CE} (identity)	D-ECE	D-ECE ₅₀
F-RCNN	52.96 \pm 0.05	75.88 \pm 0.08	34.88 \pm 0.10	17.03 \pm 0.16	28.79 \pm 0.12	35.56 \pm 0.12	16.40 \pm 0.13
F-RCNN + TS	49.25 \pm 0.07	70.93 \pm 0.09	21.43 \pm 0.03	11.52 \pm 0.18	8.85 \pm 0.08	21.43 \pm 0.03	11.37 \pm 0.23
F-RCNN + \widehat{CE} ($\lambda = 2$)	52.34 \pm 0.05	74.90 \pm 0.03	34.32 \pm 0.19	15.95 \pm 0.20	28.00 \pm 0.22	35.05 \pm 0.18	15.40 \pm 0.21
F-RCNN + \widehat{CE} ($\lambda = 3$)	51.95 \pm 0.01	74.56 \pm 0.07	34.14 \pm 0.06	15.59 \pm 0.08	27.70 \pm 0.09	34.81 \pm 0.11	14.99 \pm 0.12
F-RCNN + \widehat{CE} ($\lambda = 4$)	51.50 \pm 0.06	74.20 \pm 0.11	34.23 \pm 0.17	15.51 \pm 0.23	27.70 \pm 0.21	34.95 \pm 0.16	14.95 \pm 0.24
RetinaNet	53.20 \pm 0.02	73.36 \pm 0.02	24.11 \pm 0.09	7.95 \pm 0.19	16.97 \pm 0.07	23.92 \pm 0.08	6.92 \pm 0.15
RetinaNet + TS	53.20 \pm 0.02	73.36 \pm 0.02	24.17 \pm 0.20	7.97 \pm 0.21	17.05 \pm 0.27	24.00 \pm 0.22	6.98 \pm 0.20
RetinaNet + TCD	53.50 \pm 0.07	73.74 \pm 0.10	25.98 \pm 0.03	9.25 \pm 0.15	19.58 \pm 0.02	25.90 \pm 0.01	8.12 \pm 0.19
RetinaNet + \widehat{CE} ($\lambda = 0.1$)	52.95 \pm 0.00	72.89 \pm 0.08	23.69 \pm 0.01	7.93 \pm 0.02	16.49 \pm 0.03	23.52 \pm 0.02	7.01 \pm 0.02
RetinaNet + \widehat{CE} ($\lambda = 0.5$)	52.06 \pm 0.16	71.93 \pm 0.28	21.21 \pm 0.04	7.67 \pm 0.14	11.97 \pm 0.05	21.06 \pm 0.04	7.20 \pm 0.11
RetinaNet + \widehat{CE} ($\lambda = 1$)	50.61 \pm 0.14	70.38 \pm 0.20	19.63 \pm 0.15	11.01 \pm 0.06	7.41 \pm 0.18	19.55 \pm 0.15	10.77 \pm 0.04
FCOS	52.13 \pm 0.02	73.50 \pm 0.07	22.32 \pm 0.05	15.77 \pm 0.10	12.73 \pm 0.11	22.08 \pm 0.07	14.46 \pm 0.04
FCOS + TS	49.71 \pm 0.08	69.30 \pm 0.06	19.23 \pm 0.06	9.91 \pm 0.03	7.09 \pm 0.21	19.19 \pm 0.08	9.82 \pm 0.05
FCOS + TCD	52.26 \pm 0.03	73.58 \pm 0.07	22.55 \pm 0.16	15.18 \pm 0.19	12.61 \pm 0.22	22.21 \pm 0.18	13.65 \pm 0.16
FCOS + \widehat{CE} ($\lambda = 0.1$)	52.05 \pm 0.04	73.15 \pm 0.17	22.14 \pm 0.07	16.15 \pm 0.08	12.56 \pm 0.09	21.98 \pm 0.07	15.13 \pm 0.14
FCOS + \widehat{CE} ($\lambda = 0.5$)	51.46 \pm 0.05	72.60 \pm 0.05	22.00 \pm 0.12	16.89 \pm 0.17	12.65 \pm 0.20	21.75 \pm 0.07	15.75 \pm 0.06
FCOS + \widehat{CE} ($\lambda = 1$)	51.02 \pm 0.06	72.10 \pm 0.08	21.90 \pm 0.11	18.03 \pm 0.04	13.00 \pm 0.23	21.84 \pm 0.11	17.34 \pm 0.03

Table 9. Comparison of performance before and after finetuning on Cityscapes for 3 epochs and increasing values of λ .

Model	AP	AP ₅₀	\widehat{CE}	\widehat{CE}_{50}	\widehat{CE} (identity)	D-ECE	D-ECE ₅₀
F-RCNN	35.36 \pm 0.65	54.46 \pm 0.82	37.45 \pm 0.46	17.00 \pm 0.46	28.48 \pm 0.68	37.83 \pm 0.76	16.65 \pm 0.91
F-RCNN + \widehat{CE} ($\lambda = 1$)	35.06 \pm 0.51	53.87 \pm 0.53	35.90 \pm 0.81	15.22 \pm 0.37	25.76 \pm 0.85	35.29 \pm 1.00	14.48 \pm 0.78
F-RCNN + \widehat{CE} ($\lambda = 2$)	34.65 \pm 0.18	52.84 \pm 0.13	34.03 \pm 0.98	14.23 \pm 0.64	23.69 \pm 1.06	33.86 \pm 0.96	13.35 \pm 0.79
F-RCNN + \widehat{CE} ($\lambda = 3$)	33.90 \pm 0.30	51.52 \pm 0.49	33.99 \pm 0.79	14.80 \pm 0.32	23.19 \pm 0.98	33.39 \pm 0.90	13.33 \pm 0.62
F-RCNN + \widehat{CE} ($\lambda = 4$)	33.30 \pm 0.30	50.31 \pm 0.42	33.10 \pm 0.64	14.19 \pm 0.24	22.19 \pm 0.88	32.91 \pm 0.90	13.69 \pm 0.39
F-RCNN + \widehat{CE} ($\lambda = 5$)	33.02 \pm 0.26	50.05 \pm 0.76	32.37 \pm 0.92	14.23 \pm 0.33	20.62 \pm 1.46	32.13 \pm 1.00	13.32 \pm 0.40

Table 10. Finetuning for three epochs on COCO with several values of λ .

Model	AP	AP ₅₀	$\widehat{\text{CE}}$	$\widehat{\text{CE}}_{50}$	$\widehat{\text{CE}}$ (identity)	D-ECE	D-ECE ₅₀
F-RCNN	36.11 \pm 0.10	53.35 \pm 0.15	37.33 \pm 0.09	20.48 \pm 0.10	31.76 \pm 0.05	38.87 \pm 0.08	21.34 \pm 0.14
F-RCNN + TS	32.86 \pm 0.03	48.07 \pm 0.05	24.35 \pm 0.12	11.97 \pm 0.21	13.49 \pm 0.19	25.05 \pm 0.13	11.85 \pm 0.11
F-RCNN + $\widehat{\text{CE}}$ ($\lambda = 0.01$)	36.11 \pm 0.13	53.17 \pm 0.12	37.47 \pm 0.40	20.65 \pm 0.52	31.86 \pm 0.48	38.88 \pm 0.32	21.36 \pm 0.44
F-RCNN + $\widehat{\text{CE}}$ ($\lambda = 0.05$)	36.09 \pm 0.15	53.24 \pm 0.21	37.07 \pm 0.10	20.29 \pm 0.09	31.44 \pm 0.07	38.57 \pm 0.01	21.15 \pm 0.07
F-RCNN + $\widehat{\text{CE}}$ ($\lambda = 0.1$)	35.98 \pm 0.12	53.08 \pm 0.17	36.78 \pm 0.15	19.99 \pm 0.17	31.07 \pm 0.14	38.19 \pm 0.05	20.84 \pm 0.10
F-RCNN + $\widehat{\text{CE}}$ ($\lambda = 0.5$)	34.72 \pm 0.09	51.52 \pm 0.13	33.56 \pm 0.09	16.63 \pm 0.08	26.91 \pm 0.14	34.65 \pm 0.09	17.10 \pm 0.14
F-RCNN + $\widehat{\text{CE}}$ ($\lambda = 1$)	33.08 \pm 0.04	49.37 \pm 0.06	30.69 \pm 0.17	14.34 \pm 0.24	23.25 \pm 0.16	31.77 \pm 0.02	14.88 \pm 0.10
F-RCNN + $\widehat{\text{CE}}$ ($\lambda = 2$)	29.32 \pm 0.06	43.98 \pm 0.07	26.37 \pm 0.20	13.07 \pm 0.17	16.96 \pm 0.24	27.69 \pm 0.11	13.59 \pm 0.16
F-RCNN + $\widehat{\text{CE}}$ ($\lambda = 3$)	26.16 \pm 0.12	39.05 \pm 0.07	25.33 \pm 0.02	14.57 \pm 0.12	15.45 \pm 0.09	26.74 \pm 0.12	15.04 \pm 0.16
RetinaNet	30.83 \pm 0.12	43.33 \pm 0.22	21.89 \pm 0.36	12.03 \pm 0.19	11.03 \pm 0.56	22.24 \pm 0.27	11.63 \pm 0.08
RetinaNet + TS	30.83 \pm 0.12	43.33 \pm 0.22	25.76 \pm 0.38	10.74 \pm 0.54	17.47 \pm 0.50	26.34 \pm 0.29	10.24 \pm 0.43
RetinaNet + $\widehat{\text{CE}}$ ($\lambda = 0.01$)	31.03 \pm 0.07	43.67 \pm 0.09	21.98 \pm 0.15	11.80 \pm 0.16	11.04 \pm 0.22	22.06 \pm 0.22	10.83 \pm 0.34
RetinaNet + $\widehat{\text{CE}}$ ($\lambda = 0.05$)	31.04 \pm 0.07	43.68 \pm 0.10	21.98 \pm 0.15	11.79 \pm 0.17	11.05 \pm 0.22	22.07 \pm 0.20	10.85 \pm 0.32
RetinaNet + $\widehat{\text{CE}}$ ($\lambda = 0.1$)	31.04 \pm 0.06	43.68 \pm 0.09	21.95 \pm 0.15	11.77 \pm 0.15	11.01 \pm 0.21	22.05 \pm 0.19	10.83 \pm 0.29
RetinaNet + $\widehat{\text{CE}}$ ($\lambda = 0.5$)	30.97 \pm 0.07	43.58 \pm 0.11	21.87 \pm 0.12	11.65 \pm 0.11	10.85 \pm 0.19	22.07 \pm 0.16	10.88 \pm 0.18
RetinaNet + $\widehat{\text{CE}}$ ($\lambda = 1$)	30.77 \pm 0.05	43.30 \pm 0.05	21.79 \pm 0.16	11.94 \pm 0.03	10.63 \pm 0.33	22.09 \pm 0.21	11.32 \pm 0.07
RetinaNet + $\widehat{\text{CE}}$ ($\lambda = 2$)	30.44 \pm 0.06	42.74 \pm 0.11	21.63 \pm 0.18	12.51 \pm 0.08	10.08 \pm 0.28	21.69 \pm 0.22	11.72 \pm 0.02
RetinaNet + $\widehat{\text{CE}}$ ($\lambda = 3$)	30.10 \pm 0.06	42.23 \pm 0.10	21.58 \pm 0.14	13.14 \pm 0.16	9.72 \pm 0.18	21.48 \pm 0.06	12.29 \pm 0.32
RetinaNet + $\widehat{\text{CE}}$ ($\lambda = 4$)	29.70 \pm 0.08	41.68 \pm 0.16	21.44 \pm 0.03	13.69 \pm 0.28	9.37 \pm 0.10	21.39 \pm 0.10	12.98 \pm 0.20
FCOS	34.02 \pm 0.06	48.84 \pm 0.04	24.40 \pm 0.13	16.55 \pm 0.21	16.87 \pm 0.07	24.73 \pm 0.12	15.57 \pm 0.20
FCOS + TS	29.51 \pm 0.05	41.13 \pm 0.02	20.53 \pm 0.06	13.66 \pm 0.14	8.13 \pm 0.04	20.54 \pm 0.07	13.00 \pm 0.23
FCOS + $\widehat{\text{CE}}$ ($\lambda = 0.01$)	34.05 \pm 0.07	48.86 \pm 0.03	24.43 \pm 0.15	16.56 \pm 0.23	16.81 \pm 0.13	24.73 \pm 0.11	15.60 \pm 0.21
FCOS + $\widehat{\text{CE}}$ ($\lambda = 0.05$)	34.12 \pm 0.03	48.89 \pm 0.01	24.29 \pm 0.16	16.29 \pm 0.13	16.68 \pm 0.14	24.65 \pm 0.14	15.37 \pm 0.19
FCOS + $\widehat{\text{CE}}$ ($\lambda = 0.1$)	34.03 \pm 0.08	48.83 \pm 0.04	24.18 \pm 0.13	16.34 \pm 0.12	16.56 \pm 0.13	24.57 \pm 0.11	15.40 \pm 0.14
FCOS + $\widehat{\text{CE}}$ ($\lambda = 0.5$)	33.89 \pm 0.09	48.65 \pm 0.08	24.09 \pm 0.12	16.35 \pm 0.09	16.46 \pm 0.10	24.45 \pm 0.09	15.39 \pm 0.05
FCOS + $\widehat{\text{CE}}$ ($\lambda = 1$)	33.73 \pm 0.12	48.45 \pm 0.11	23.95 \pm 0.12	16.35 \pm 0.10	16.25 \pm 0.08	24.27 \pm 0.14	15.36 \pm 0.06
FCOS + $\widehat{\text{CE}}$ ($\lambda = 2$)	33.41 \pm 0.13	48.00 \pm 0.15	23.41 \pm 0.22	16.15 \pm 0.25	15.56 \pm 0.16	23.83 \pm 0.12	15.27 \pm 0.23
FCOS + $\widehat{\text{CE}}$ ($\lambda = 3$)	33.02 \pm 0.12	47.44 \pm 0.14	23.18 \pm 0.19	16.34 \pm 0.17	15.32 \pm 0.17	23.72 \pm 0.17	15.67 \pm 0.21



Figure 5. Qualitative results on COCO using uncalibrated (left) and fine-tuned with \widehat{CE} (right) F-RCNN model. Detection threshold is set to 0.5. Red boxes denote predictions with their corresponding confidence scores. Ground truth boxes are shown in green dashed lines. The calibration loss is used with identity link, i.e., the scores should be aligned with the IoU overlap with a ground truth box. Best viewed in color and zoom.