# Supplementary Material

**Table S1.** Maximum Calibration Error

|  | ET | TC | WT | Avg |
|---|---|---|---|---|
| CE | $0.26 \pm 0.19$ | $0.28 \pm 0.18$ | $0.23 \pm 0.12$ | $0.25 \pm 0.12$ |
| CE + Ts | $0.30 \pm 0.16$ | $0.29 \pm 0.16$ | $0.25 \pm 0.10$ | $0.28 \pm 0.10$ |
| CE + L1-ACE | $0.24 \pm 0.19$ | $0.26 \pm 0.18$ | $0.20 \pm 0.15$ | $0.23 \pm 0.14$ |
| CE + L1-ACE + Ts | $\mathbf{0.23 \pm 0.15}$ | $\mathbf{0.25 \pm 0.16}$ | $0.20 \pm 0.12$ | $\mathbf{0.23 \pm 0.11}$ |
| Dice | $0.59 \pm 0.15$ | $0.61 \pm 0.15$ | $0.55 \pm 0.13$ | $0.58 \pm 0.11$ |
| Dice + Ts | $0.56 \pm 0.15$ | $0.58 \pm 0.15$ | $0.51 \pm 0.14$ | $0.55 \pm 0.11$ |
| Dice + L1-ACE | $0.29 \pm 0.18$ | $0.30 \pm 0.17$ | $0.23 \pm 0.14$ | $0.27 \pm 0.13$ |
| Dice + L1-ACE + Ts | $\mathbf{0.23 \pm 0.17}$ | $\mathbf{0.25 \pm 0.16}$ | $0.21 \pm 0.11$ | $\mathbf{0.23 \pm 0.12}$ |
| Dice-CE | $0.46 \pm 0.17$ | $0.46 \pm 0.17$ | $0.36 \pm 0.17$ | $0.43 \pm 0.13$ |
| Dice-CE + Ts | $0.26 \pm 0.19$ | $0.32 \pm 0.19$ | $\mathbf{0.19 \pm 0.12}$ | $0.26 \pm 0.13$ |
| Dice-CE + L1-ACE | $0.29 \pm 0.21$ | $0.30 \pm 0.20$ | $0.21 \pm 0.14$ | $0.27 \pm 0.15$ |
| Dice-CE + L1-ACE + Ts | $0.25 \pm 0.19$ | $0.27 \pm 0.18$ | $0.20 \pm 0.11$ | $0.24 \pm 0.13$ |

**Table S2.** Expected Calibration Error ($\times 10^{-4}$)

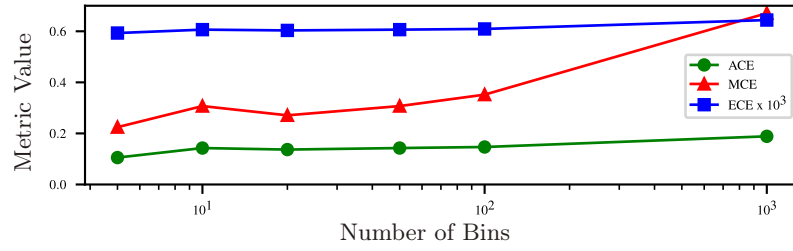|  | ET | TC | WT | Avg |
|---|---|---|---|---|
| CE | $4.21 \pm 7.49$ | $5.66 \pm 10.5$ | $\mathbf{9.71 \pm 11.3}$ | $6.53 \pm 7.53$ |
| CE + Ts | $7.77 \pm 7.16$ | $10.2 \pm 10.0$ | $19.1 \pm 9.45$ | $12.4 \pm 6.93$ |
| CE + L1-ACE | $4.00 \pm 4.75$ | $6.16 \pm 9.41$ | $12.0 \pm 12.9$ | $7.38 \pm 6.94$ |
| CE + L1-ACE + Ts | $21.1 \pm 4.69$ | $27.2 \pm 6.65$ | $62.0 \pm 11.1$ | $36.8 \pm 5.69$ |
| Dice | $5.09 \pm 7.67$ | $7.48 \pm 16.5$ | $15.5 \pm 17.8$ | $9.34 \pm 11.1$ |
| Dice + Ts | $4.95 \pm 7.63$ | $7.36 \pm 16.5$ | $15.0 \pm 17.8$ | $9.10 \pm 11.1$ |
| Dice + L1-ACE | $3.47 \pm 7.56$ | $4.50 \pm 8.45$ | $10.1 \pm 12.8$ | $\mathbf{6.03 \pm 6.84}$ |
| Dice + L1-ACE + Ts | $\mathbf{3.27 \pm 7.17}$ | $\mathbf{4.41 \pm 8.00}$ | $12.0 \pm 10.9$ | $6.54 \pm 5.84$ |
| Dice-CE | $4.40 \pm 7.57$ | $5.70 \pm 12.2$ | $11.6 \pm 14.4$ | $7.21 \pm 8.62$ |
| Dice-CE + Ts | $3.84 \pm 7.43$ | $5.20 \pm 11.8$ | $10.2 \pm 13.4$ | $6.41 \pm 8.11$ |
| Dice-CE + L1-ACE | $3.91 \pm 8.09$ | $6.04 \pm 14.6$ | $10.3 \pm 14.8$ | $6.75 \pm 9.84$ |
| Dice-CE + L1-ACE + Ts | $4.23 \pm 7.68$ | $6.56 \pm 14.0$ | $13.9 \pm 12.6$ | $8.23 \pm 8.83$ |

**Fig. S1.** Comparison of calibration metrics using different number of bins when evaluated on the DSC + mL1-ACE loss model. Both ACE and ECE are very stable with respect to number of bins (5, 10, 20, 50, 100, 1000). MCE shows a larger increase from 100 to 1000 bins, this is due to the presence of empty bins, a likely occurrence when using such a high number of bins.