

On Focal Loss for Class-Posterior Probability Estimation: A Theoretical Perspective

Nontawat Charoenphakdee*^{1,2}, Jayakorn Vongkulbhisal*³, Nuttapong Chairatanakul^{4,5}, Masashi Sugiyama^{2,1}
 The University of Tokyo¹, RIKEN AIP², IBM Research³, Tokyo Institute of Technology⁴, RWBC-OIL (AIST)⁵



Summary

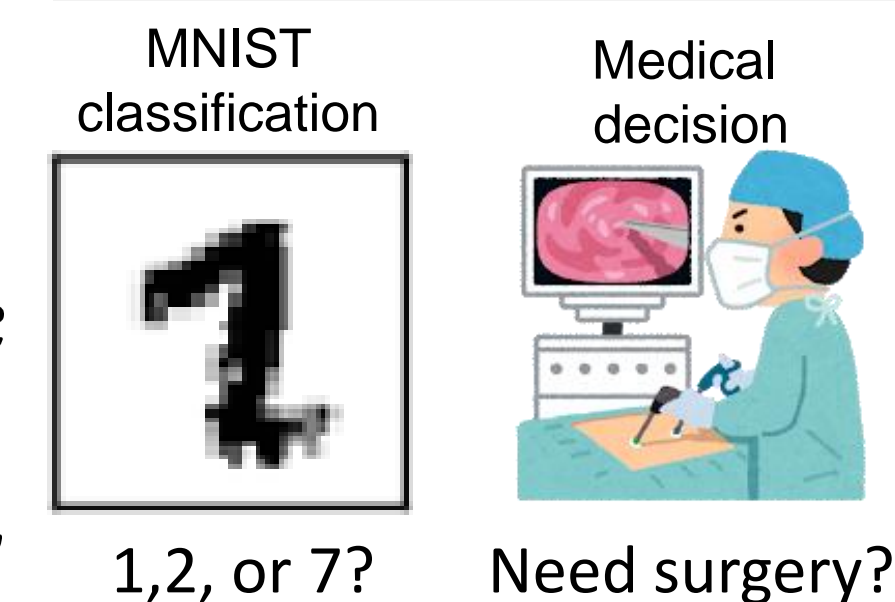
Theoretical analysis of focal loss with practical use.

Q1: Does focal risk minimizer give Bayes-optimal classifier?
Yes!

Q2: Does focal risk minimizer match class-posterior probability $p(y|x)$?
No! Directly using model's output gives **unreliable confidence**.

Q3: Following Q2, can we do anything about it?
Yes! We discovered a closed-form transformation Ψ^γ that can recover $p(y|x)$ **with theoretical guarantee!**

Introduction



	$p(y x)$	$q_{CE}^*(x)$	$q_{FL,1}^*(x)$	$q_{FL,3}^*(x)$	$q_{FL,5}^*(x)$
$y = +1$	$\begin{pmatrix} .90 \\ .10 \end{pmatrix}$	$\begin{pmatrix} .90 \\ .10 \end{pmatrix}$	$\begin{pmatrix} .78 \\ .22 \end{pmatrix}$	$\begin{pmatrix} .65 \\ .35 \end{pmatrix}$	$\begin{pmatrix} .60 \\ .40 \end{pmatrix}$
$y = -1$	$\begin{pmatrix} .10 \\ .90 \end{pmatrix}$	$\begin{pmatrix} .10 \\ .90 \end{pmatrix}$	$\begin{pmatrix} .22 \\ .78 \end{pmatrix}$	$\begin{pmatrix} .35 \\ .65 \end{pmatrix}$	$\begin{pmatrix} .40 \\ .60 \end{pmatrix}$

$q_\ell^* = \operatorname{argmin}_q \mathbb{E}_{y \sim p(y|x)} [\ell(q(x), e_y)]$
 q_{CE}^* : Cross-entropy risk minimizer
 $q_{FL,\gamma}^*$: Focal risk minimizer

$q(x) \in \Delta^K$
 e_y : One-hot vector

- Bayes-optimal classifier predicts the most probable class $\arg \max_y p(y|x)$.
- Class-posterior probability $p(y|x)$ provides useful confidence score.
- **Loss function highly influences the behavior of the trained model.**

Example: The well-studied cross-entropy (CE) loss for K-class classification:

$$\ell_{CE}(v, u) = - \sum_{i=1}^K u_i \log(v_i). \quad u \in \Delta^K, v \in \Delta^K$$

u is often a one-hot vector

CE loss is

- **classification-calibrated:** CE risk minimizer $q_{CE}^*(x)$ gives Bayes optimal classifier.
- **strictly proper:** CE risk minimizer $q_{CE}^*(x)$ gives class-posterior probability.

Q: What about theoretical properties of focal loss?

Focal loss

$$\ell_{FL}^\gamma(v, u) = - \sum_{i=1}^K u_i (1 - v_i)^\gamma \log(v_i) \quad \gamma \geq 0$$

(Lin+, 2017)

- Originally proposed for dense object detection.
- Many practical applications in the medical field.

(Al Rahhal+, 2019, Chang+, 2018, Lotfy+, 2019, Sun+, 2019)

Main results

Focal loss is **classification-calibrated** for $\gamma \geq 0$:

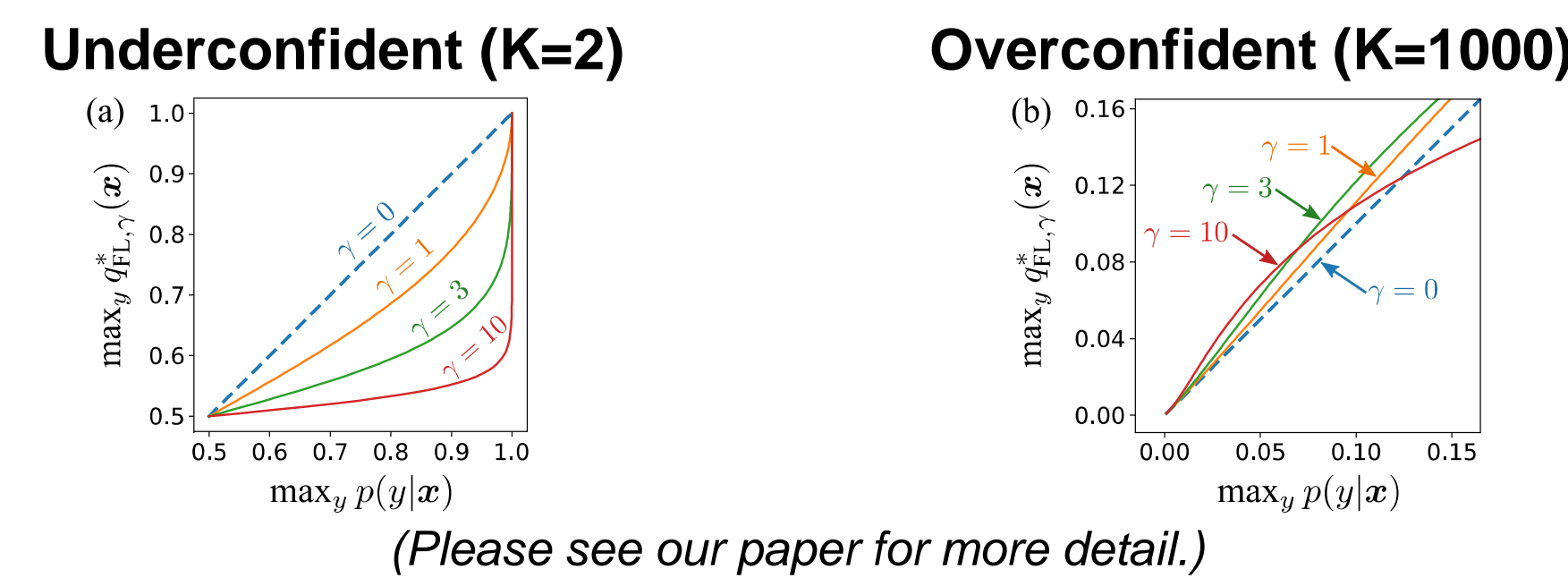
$$\arg \max_y q_{FL,\gamma}^*(x) = \arg \max_y p(y|x).$$

However, it is **not strictly proper** for $\gamma > 0$:

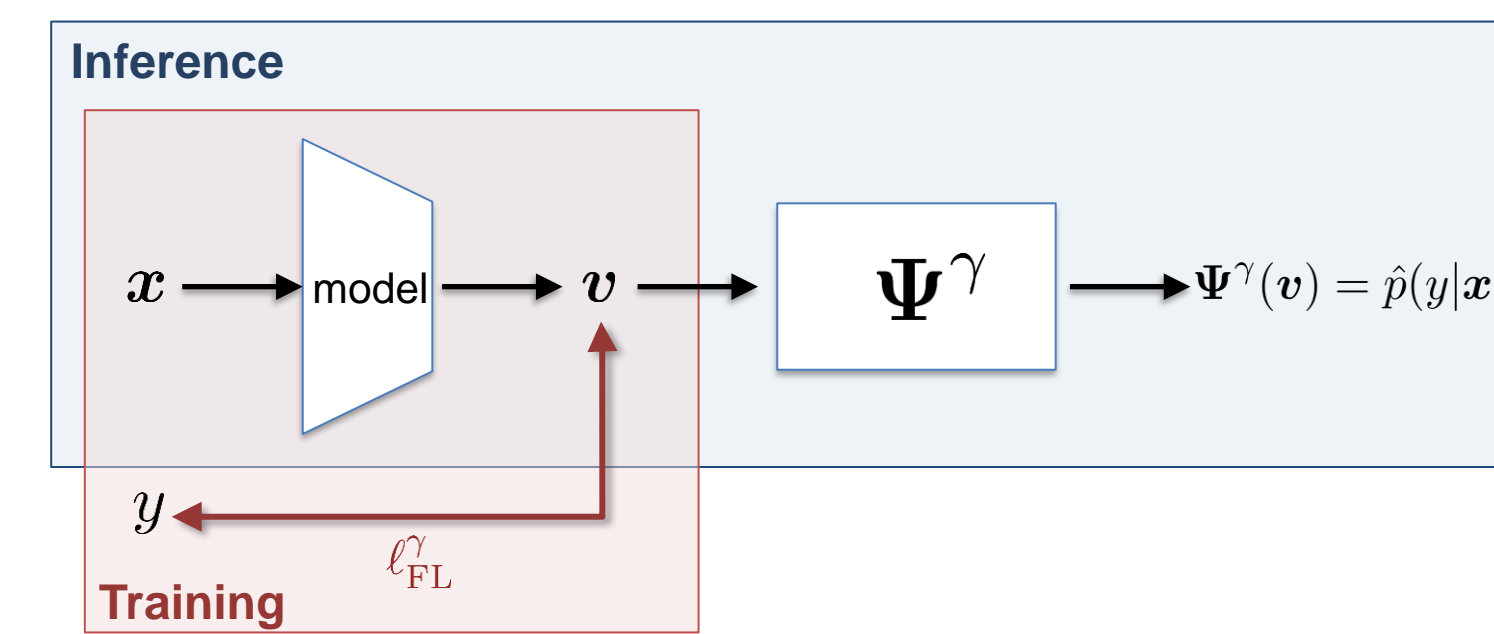
$$q_{FL,\gamma}^*(x) \neq p(y|x).$$

We can predict the most probable class, but **confidence score is unreliable**.

Example:



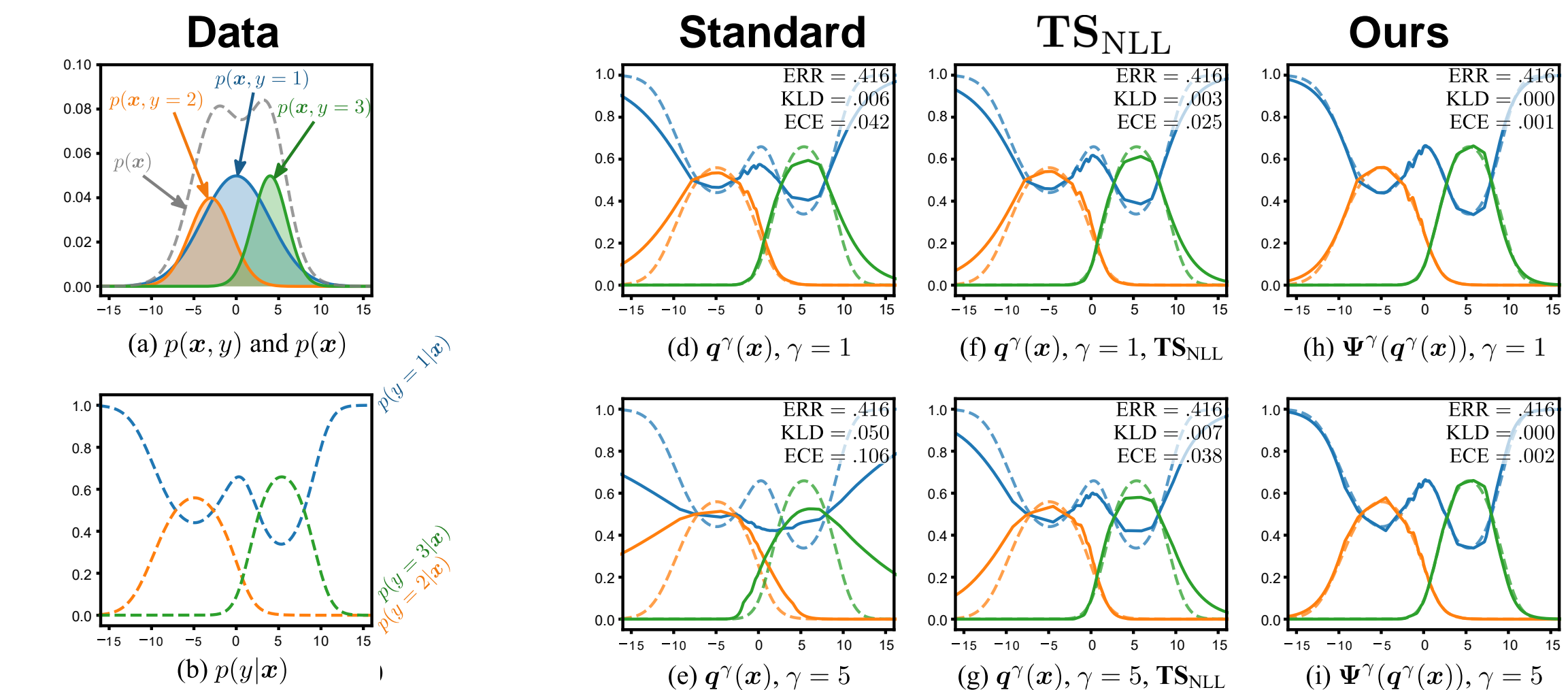
Solution: Recover $p(y|x)$ from $q_{FL,\gamma}^*(x)$ via Ψ^γ :



Define $\Psi^\gamma(v) = [\Psi_1^\gamma(v), \dots, \Psi_K^\gamma(v)]^\top$,
 where $\Psi_i^\gamma(v) = \frac{h^\gamma(v_i)}{\sum_{i=1}^K h^\gamma(v_i)}$,
 and $h^\gamma(v_i) = \frac{v_i}{(1-v_i)^\gamma - \gamma(1-v_i)^{\gamma-1} v_i \log v_i}$.

- Closed-form
- No hyperparameter
- Theoretically justified
- Preserves accuracy
- No additional training required

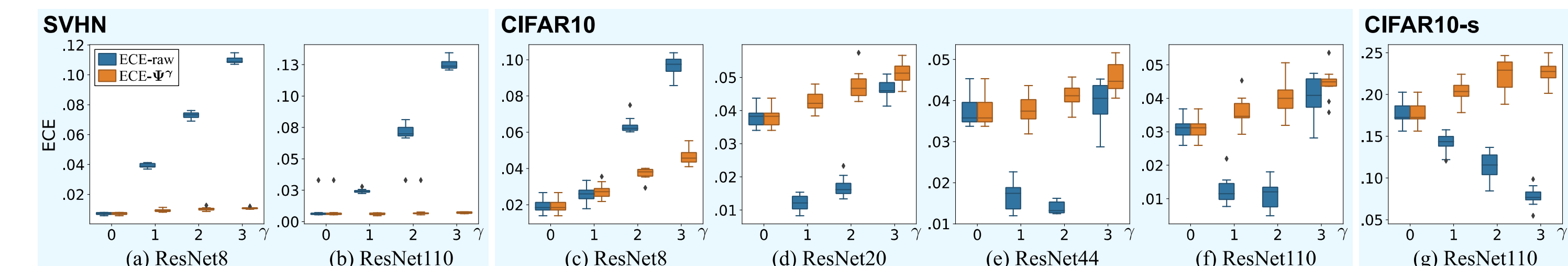
Numerical simulation



- Bigger γ makes the network q^γ more prone to be underconfident in **Standard**.
- Using temperature scaling TS_{NLL} (Guo+, 2017) is insufficient to recover $p(y|x)$.
- With our proposed Ψ^γ , we can recover $p(y|x)$ (almost) perfectly.

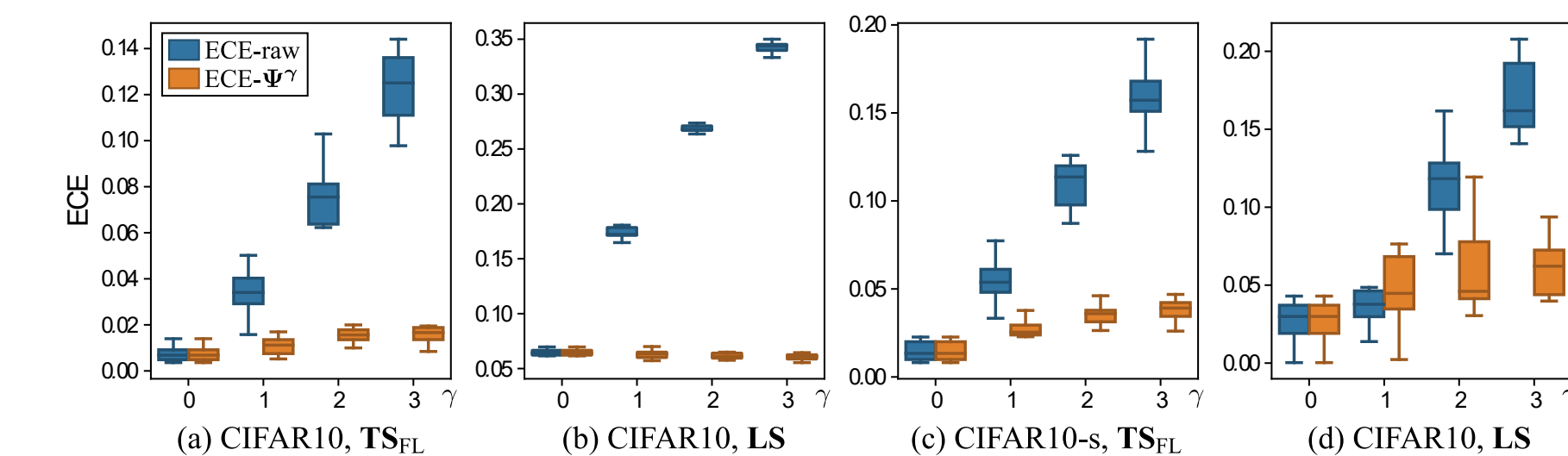
Benchmark experiments

Evaluation metric: Expected calibration error (ECE) (Naeini+, 2015)



- Using Ψ^γ is effective when we have good approximation of $q_{FL,\gamma}^*(x)$ (Fig. a-c)
- But it is less effective when having small data or model architecture is too large (Fig. d-g)

With focal-loss-based temperature scaling TS_{FL} or label smoothing LS :



Using Ψ^γ is preferable for both cases.

*We used ResNet110 for Fig. a-d. Same trend can be observed for all models in our paper (ResNet8-110).

References

Guo, C., Pleiss, G., Sun, Y., and Weinberger, K. Q. On calibration of modern neural networks. In ICML, pages 1321–1330. JMLR, 2017.
 Lin, T. Y., Goyal, P., Girshick, R., He, K., & Dollár, P. Focal loss for dense object detection. ICCV, 2017.
 Naeini, M. P., Cooper, G., and Hauskrecht, M. Obtaining well-calibrated probabilities using bayesian binning. AAAI, 2015.
 Al Rahhal, M. M., Bazi, Y., Almubarak, H., Alajlan, N. and Al Zuair, M. Dense convolutional networks with focal loss and image generation for electrocardiogram classification. IEEE Access, 2019.
 Chang, J., Zhang, X., Ye, M., Huang, D., Wang, P. and Yao, C. Brain tumor segmentation based on 3D Unet with multi-class focal loss. CISP-BMEI, 2018.
 Lotfy, M., Shubair, R. M., Navab, N. and Albarqouni, S. Investigation of Focal Loss in Deep Learning Models For Femur Fractures Classification. ICCTA, 2019.
 Sun, X., Dong, K., Ma, L., Sutcliffe, R., He, F., Chen, S. and Feng, J. Drug-drug interaction extraction via recurrent hybrid convolutional neural networks with an improved focal loss. Entropy, 2019.