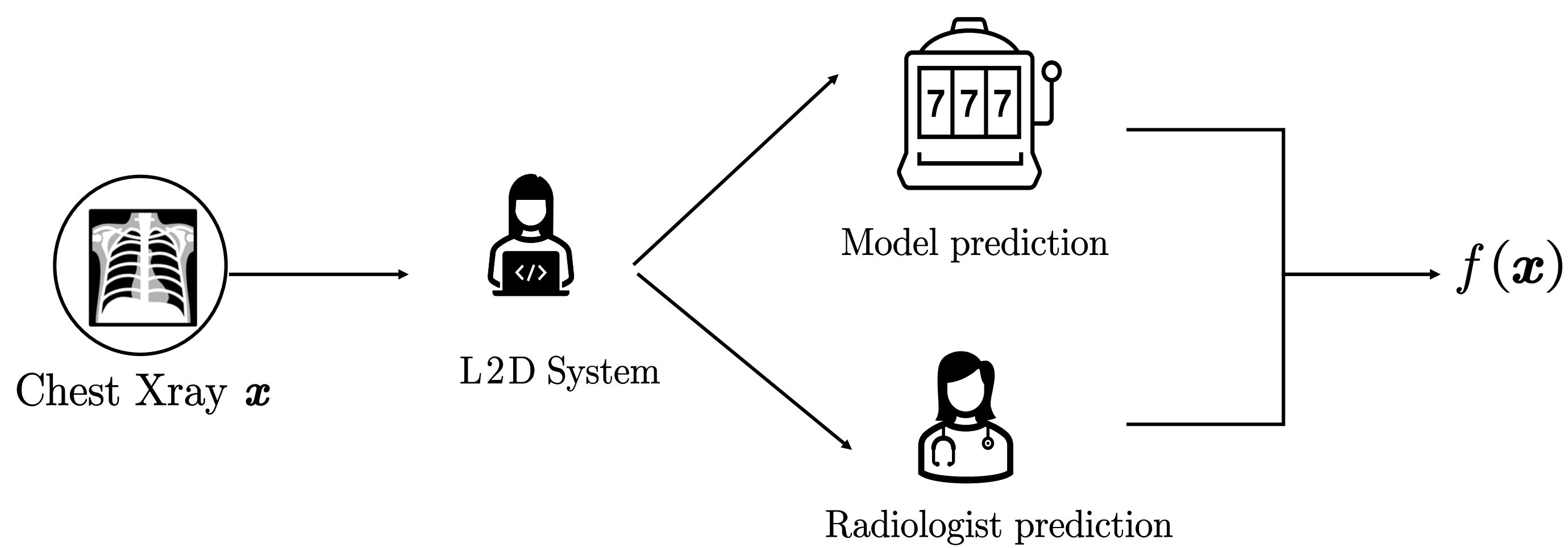


## Learning to Defer (L2D)

Learning to defer (L2D) allows the classifier to defer its prediction to an expert for safer predictions, by balancing the system's accuracy and extra costs incurred by consulting the expert.



$$\text{Find } f^* = \underset{f}{\operatorname{argmin}} \mathbb{E}_{p(\mathbf{x}, y, m)} [\mathbb{I}(f(\mathbf{x}) \neq \perp \text{ and } f(\mathbf{x}) \neq y) + \mathbb{I}(f(\mathbf{x}) = \perp \text{ and } m \neq y)] + \mathbb{E}_{p(\mathbf{x})} [c\mathbb{I}(f(\mathbf{x}) = \perp)] \Rightarrow \text{Extra deferral cost} \quad \text{Classification Error}$$

### Goal:

- Optimize Classifier to adapt to Human's weaknesses and strengths.

### Intuitive Explanation of $f^*$ :

- $f^*(x)$  is  $\perp$ , i.e., choose to defer to expert and output  $m$ , if expert's accuracy against classifier can offset the cost:

$$\text{Acc}(m|\mathbf{x}) = \Pr(m = y|\mathbf{x}) \geq \max_y p(y|\mathbf{x}) + c$$

- Otherwise  $f^*(x)$  is the Bayes optimal  $y^* = \operatorname{argmax}_y p(y|\mathbf{x})$

## Underfitting Issues in L2D

### Existing L2D Surrogates

- Consistent loss formulation [5]

$$\min_g \mathbb{E}_{p(\mathbf{x}, y, m)} \left[ \Psi(\mathbf{g}(\mathbf{x}), y) + c\mathbb{I}(m \neq y) \sum_{y' \in [K]} \Psi(\mathbf{g}(\mathbf{x}), y') + (1-c)\mathbb{I}(m = y)\Psi(\mathbf{g}(\mathbf{x}), K+1) \right]$$

- Treat deferral option  $\perp$  as a new class
- $\Psi$  is any consistent multiclass loss
- $\mathbf{g}$  is a  $K+1$  dimensional scoring function. First  $K$  dimensions  $\mathbf{g}_{[1:K]}$  is a classifier.
- Essence: Problem reduction to  $K+1$  class classification on distribution  $\hat{p}$  [4]

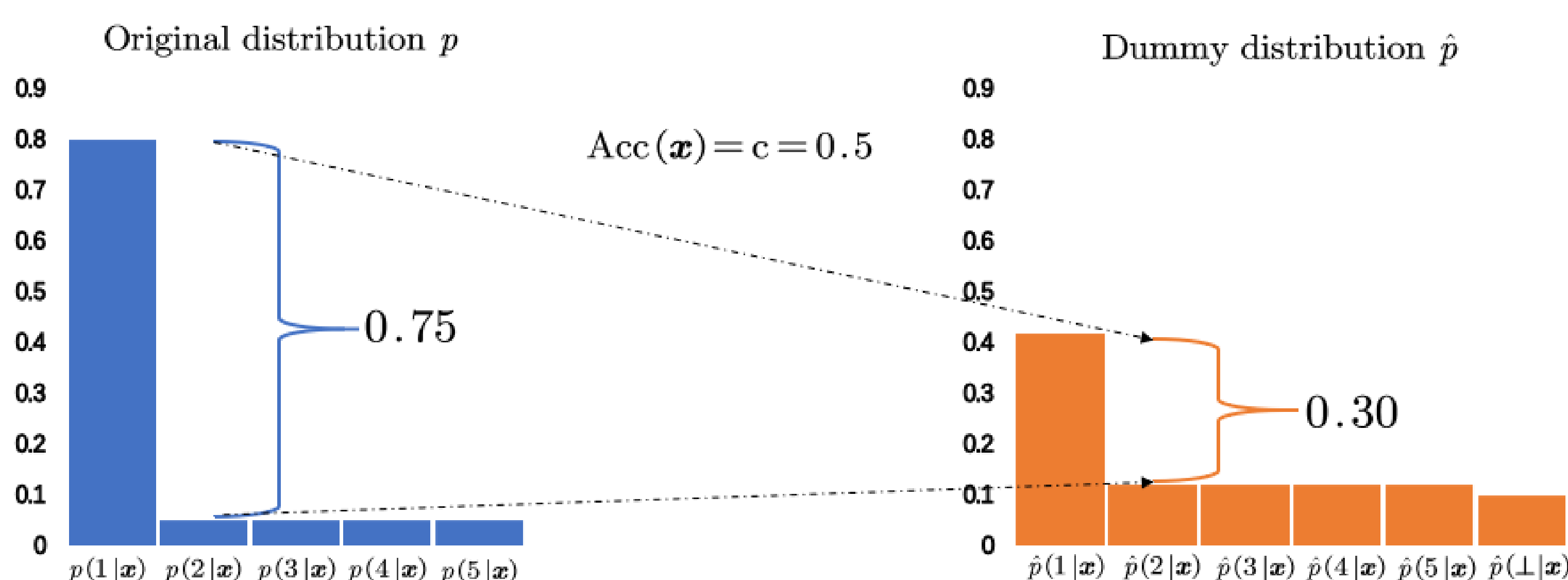
$$\{\hat{p}(y|\mathbf{x})\}_{y \in [K+1]} = \left\{ \frac{p(y|\mathbf{x}) + c(1 - \text{Acc}(m|\mathbf{x}))}{\hat{\xi}(\mathbf{x})} \right\}_{y \in [K]} \cup \left\{ \frac{(1-c)\text{Acc}(m|\mathbf{x})}{\hat{\xi}(\mathbf{x})} \right\}$$

### Underfitting of Classifier [3]

Sample efficiency and classifier's performance  $\mathbf{g}_{[1:K]}$  degrade rapidly with increasing extra expert cost  $c$ , i.e., underfitting occurs.

### Cause of underfitting: Probability margin shrinks

$$p(y^*|\mathbf{x}) - p(y|\mathbf{x}) > \frac{p(y^*|\mathbf{x}) - p(y|\mathbf{x})}{1 + (1-c)\text{Acc}(\mathbf{x}) + Kc(1 - \text{Acc}(\mathbf{x}))} = \hat{p}(y^*|\mathbf{x}) - \hat{p}(y|\mathbf{x})$$



$\hat{p}$  is more ambiguous!

- Direct fix:** Post-hoc methods train base model with  $c=0$ , then retrain/reformulate deferral rule to account for non-zero  $c$  [3]

## Novel Surrogate Loss

### Motivation: Eliminate redundancy in previous problem reduction

- Previous deferral rule: compare each label's acc to expert acc  
 $\text{Acc}(m|\mathbf{x}) \geq p(1|\mathbf{x}) + c, \text{Acc}(m|\mathbf{x}) \geq p(2|\mathbf{x}) + c, \dots, \text{Acc}(m|\mathbf{x}) \geq p(K|\mathbf{x}) + c,$
- Ideal deferral rule: only compare the Bayes optimal one

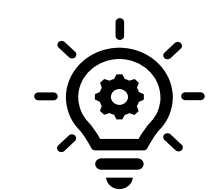
$$\text{Acc}(m|\mathbf{x}) \geq p(y^*|\mathbf{x}) + c \quad \text{Bayes optimal label}$$

- Induced loss formulation:

$$\min_g \mathbb{E}_{p(\mathbf{x}, y, m)} [\Psi(\mathbf{g}(\mathbf{x}), y) + c\mathbb{I}(m \neq y)\Psi(\mathbf{g}(\mathbf{x}), y^*) + (1-c)\mathbb{I}(m = y)\Psi(\mathbf{g}(\mathbf{x}), K+1)]$$

- Essence: Problem reduction to  $K+1$  class classification on distribution  $\tilde{p}$ :

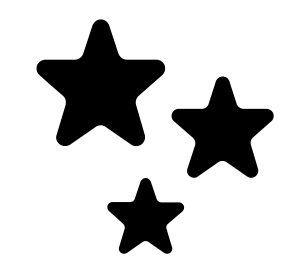
$$\{\tilde{p}(y|\mathbf{x})\}_{y \in [K]} = \left\{ \frac{p(y|\mathbf{x}) + c(1 - \text{Acc}(m|\mathbf{x}))\mathbb{I}(y = y^*)}{\tilde{\xi}(\mathbf{x})} \right\}_{y \in [K]} \cup \left\{ \frac{(1-c)\text{Acc}(m|\mathbf{x})}{\tilde{\xi}(\mathbf{x})} \right\}$$



### Result 1: Better probability margin!

$$\tilde{p}(y^*|\mathbf{x}) - \tilde{p}(y|\mathbf{x}) > \hat{p}(y^*|\mathbf{x}) - \hat{p}(y|\mathbf{x})$$

However, there is no access to the Bayes optimal label



### Approximate Bayes optimal with intermediate learning results $y'$

Def. 1: Redundancy-free Loss Formulation  $R_{\Psi}^{\perp}(\mathbf{g})$

$$\mathbb{E}_{p(\mathbf{x}, y, m)} [\Psi(\mathbf{g}(\mathbf{x}), y) + c\mathbb{I}(m \neq y) \min_{y' \in [K]} \Psi(\mathbf{g}(\mathbf{x}), y') + (1-c)\mathbb{I}(m = y)\Psi(\mathbf{g}(\mathbf{x}), K+1)]$$

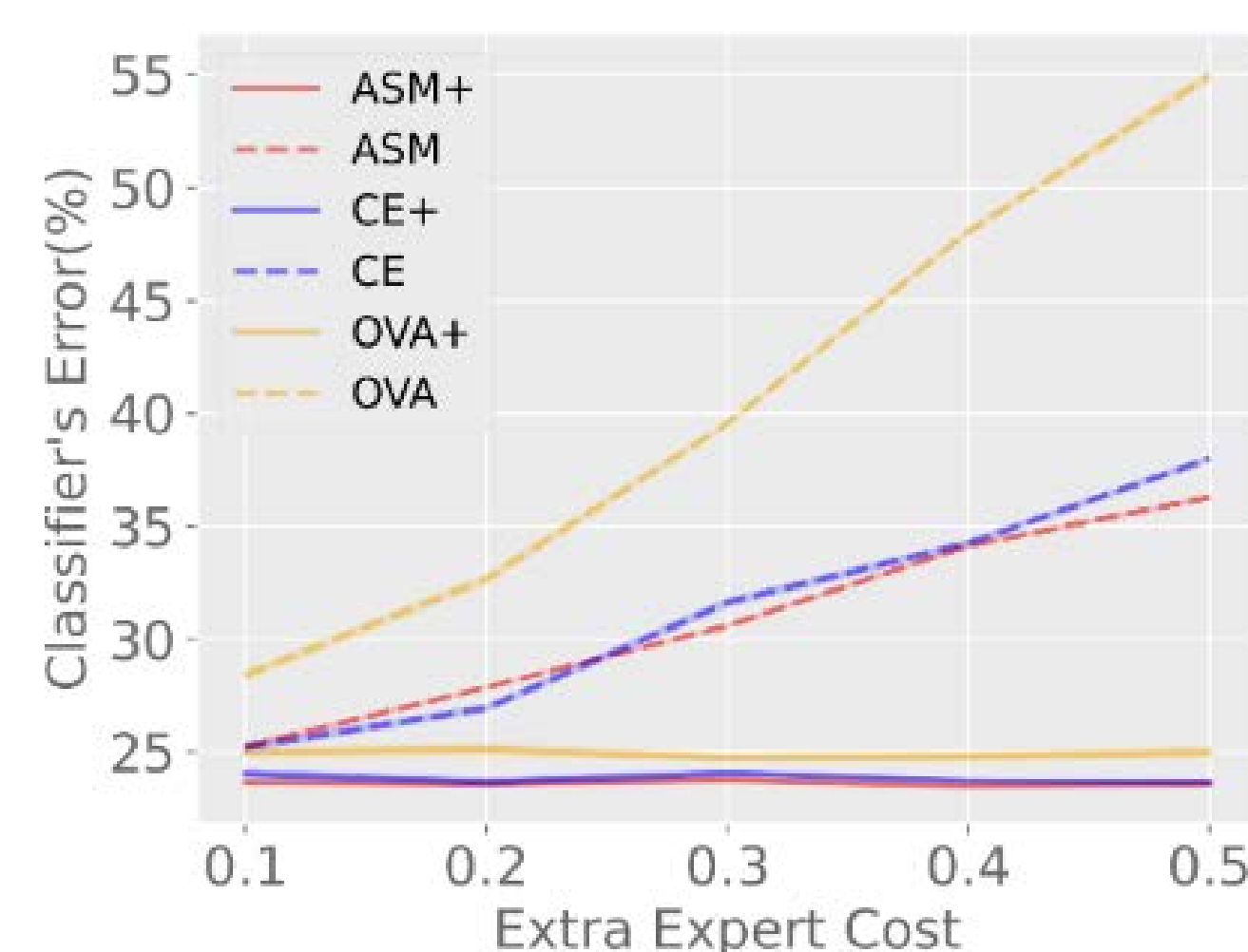


### Result 2: Consistency analysis for $\Psi$ in [1,2,6]

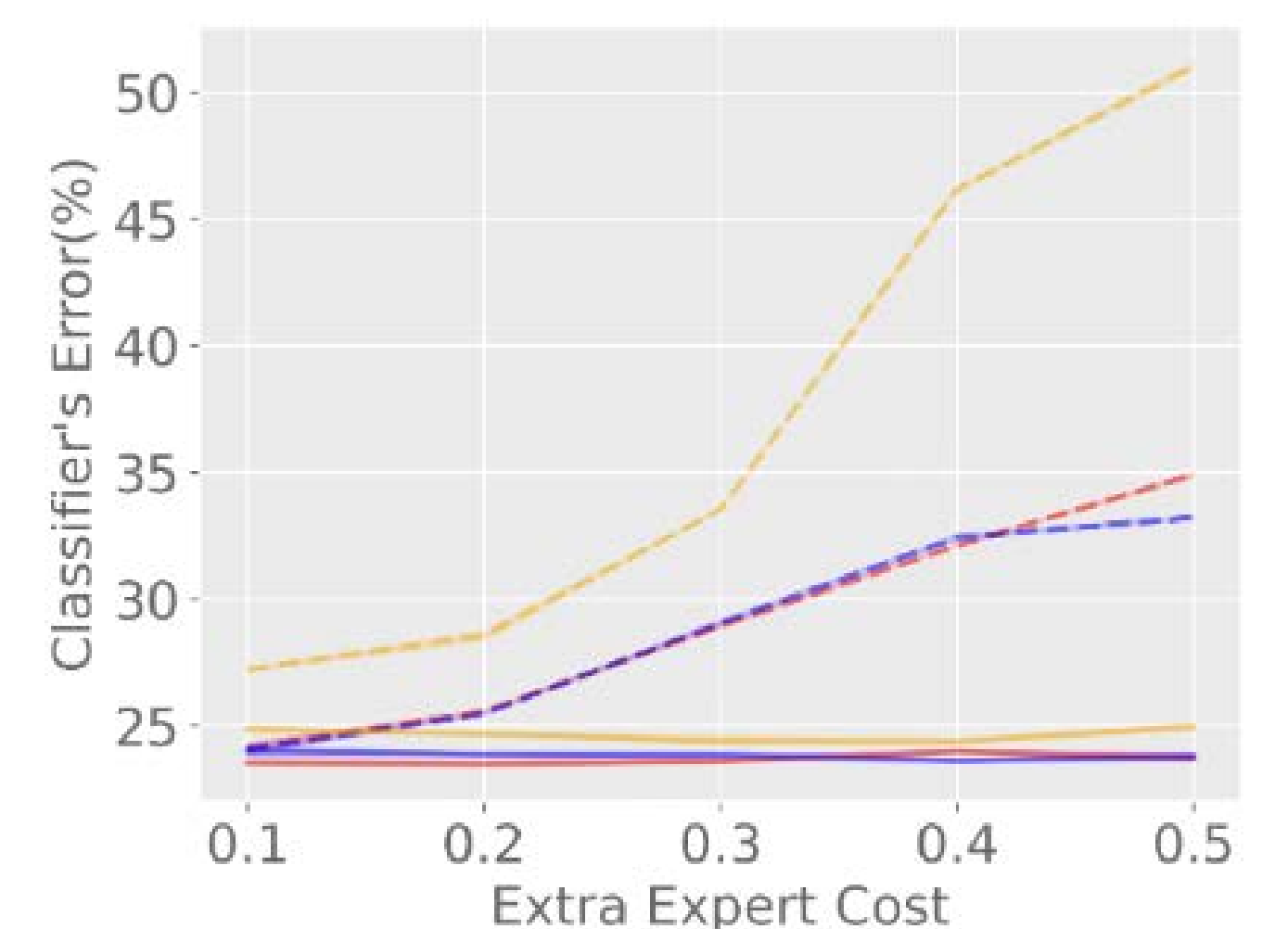
$$\text{Thm. 1: } \forall \mathbf{g}^* \in \operatorname{argmin}_{\mathbf{g}} R_{\Psi}^{\perp}(\mathbf{g}), f_{\mathbf{g}^*} = f^*$$

## Experiments

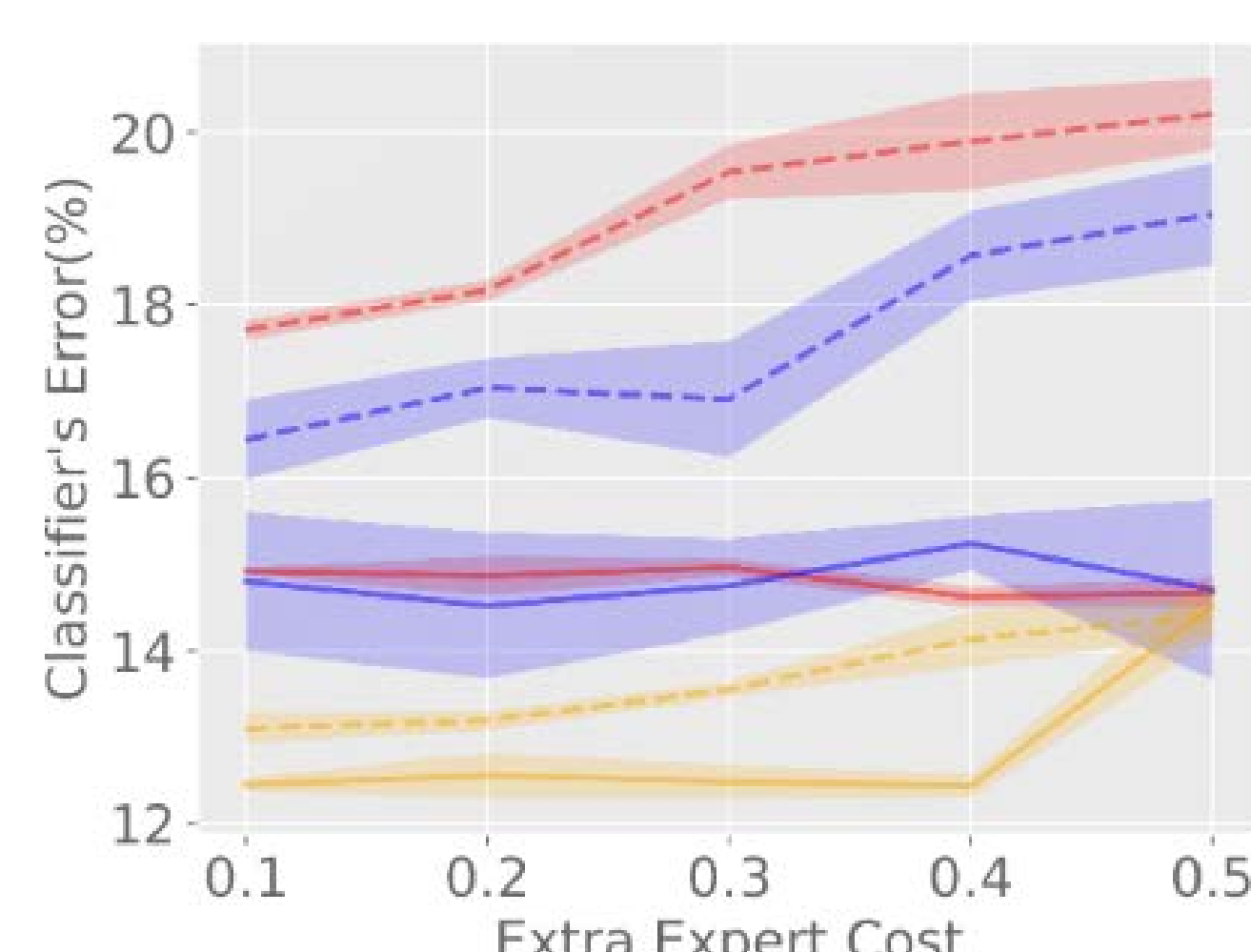
Proposed loss mitigates underfitting and generates more accurate predictions!



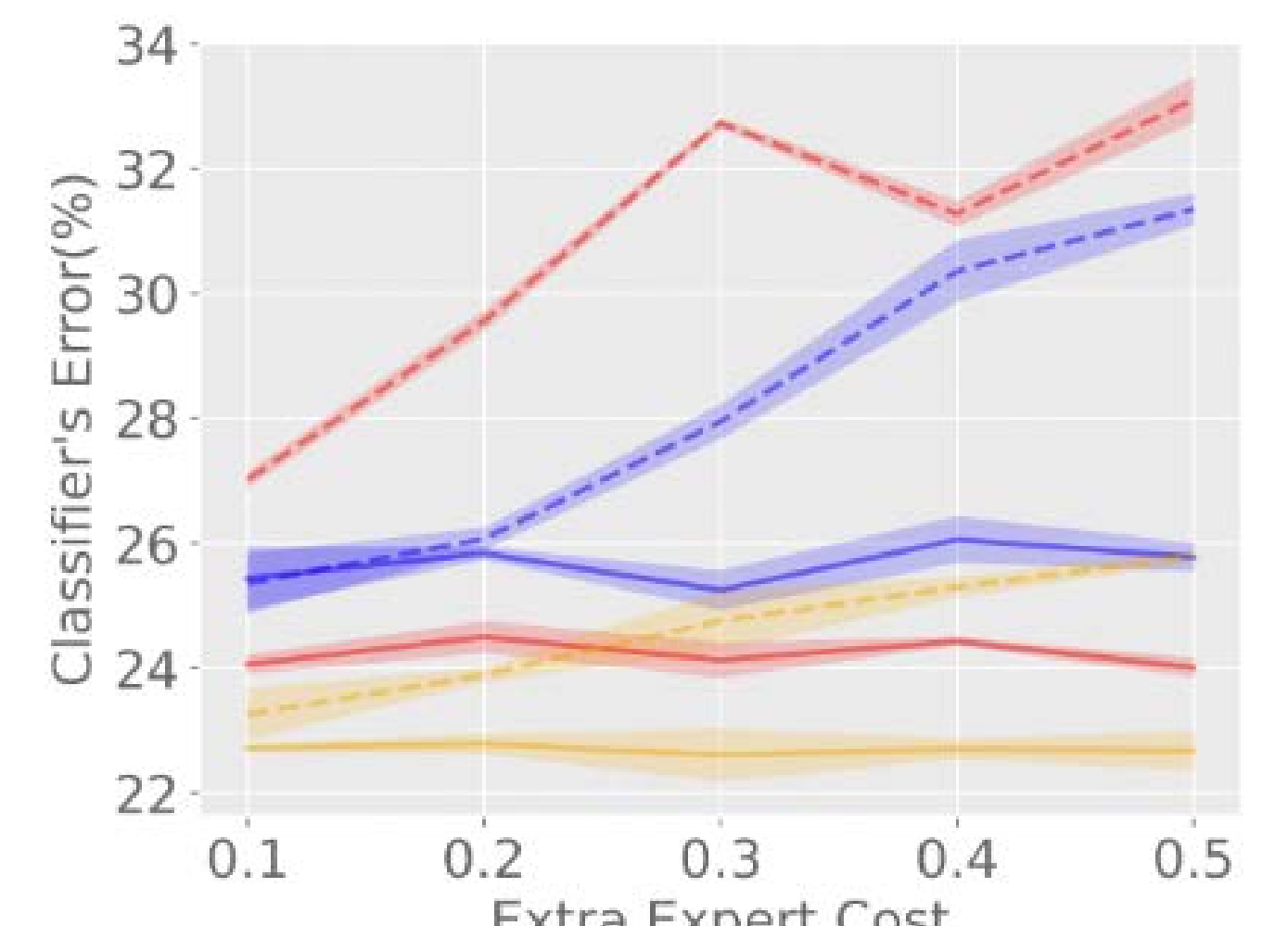
(a) CIFAR-100 (94%)



(b) CIFAR-100 (75%)



(c) CIFAR-10 (w)



(d) CIFAR-10 (o)

## References

- [1] Mozannar, H. and Sontag, D. A. (2020). Consistent estimators for learning to defer to an expert. In ICML.
- [2] Verma, R. and Nalisnick, E. T. (2022). Calibrated learning to defer with one-vs-all classifiers. In ICML.
- [3] Narasimhan, H., Jitkrittum, W., Menon, A. K., Rawat, A. S., and Kumar, S. (2022). Post-hoc estimators for learning to defer to an expert. In NeurIPS.
- [4] Cao, Y., Cai, T., Feng, L., Gu, L., GU, J., An, B., Niu, G., and Sugiyama, M. (2022). Generalizing consistent multi-class classification with rejection to be compatible with arbitrary losses. In NeurIPS.
- [5] Charusaie, M., Mozannar, H., Sontag, D. A., and Samadi, S. (2022). Sample efficient learning of predictors that complement humans. In ICML.
- [6] Cao, Y., Mozannar, H., Feng, L., Wei, H., and An, B. (2023). In defense of softmax parametrization for calibrated and consistent learning to defer. In NeurIPS.

## Link

QR code for the paper

