

## 454 A Proof of Theorem 1

455 We first prove that if there is an L2D probability estimator  $\hat{\eta}$  for  $L_\phi$ , then we can reconstruct a  
 456 probability estimator  $\tilde{\eta}$  for multiclass classification with  $\phi$  if  $\eta$  is in  $\tilde{\Delta}^{K+1}$ , which is the collection of  
 457 all the elements  $\beta \in \Delta^{K+1}$  with  $\beta_{K+1} \leq \frac{1}{2}$ . Then we show that according to the symmetry of  $\phi$ , we  
 458 can extend this estimator to  $\Delta^{K+1}$ , and then construct an unbounded L2D probability estimator for  
 459  $L_\phi$ .

460 *Proof.* Denote by  $R_{L_\phi}(\mathbf{u}, \boldsymbol{\eta}(\mathbf{x}), \Pr(M = Y|X = \mathbf{x})) = \sum_{y=1}^K \eta(\mathbf{x})_y \phi(\mathbf{u}, y) + \Pr(M = Y|X =$   
 461  $\mathbf{x})\phi(\mathbf{u}, K+1)$  the conditional risk for L2D. Denote by  $\hat{\eta}$  a probability estimator for  $L_\phi$  that  $\hat{\eta}(\mathbf{u}^*) =$   
 462  $[\boldsymbol{\eta}(\mathbf{x}); \Pr(M = Y|X = \mathbf{x})]$  for all  $\mathbf{u}^* \in \operatorname{argmin}_{\mathbf{u}} R_{L_\phi}(\mathbf{u}, \boldsymbol{\eta}(\mathbf{x}), \Pr(M = Y|X = \mathbf{x}))$ , we can  
 463 learn that  $\Delta^K \times [0, 1]$  is in the range of  $\hat{\eta}$ . Then it is easy to verify that  $\tilde{\eta}(\mathbf{u}) = \frac{\hat{\eta}(\mathbf{u})}{1 + \hat{\eta}_{K+1}(\mathbf{u})} \in \tilde{\Delta}^{K+1}$   
 464 is a valid probability estimator for  $K + 1$ -class multiclass classification with  $\phi$  if the posterior  
 465 probabilities  $[p(y|\mathbf{x})]_{y=1}^{K+1}$  is in  $\tilde{\Delta}^{K+1}$ .

Then we construct a valid estimator for  $\Delta^{K+1}$  based on  $\tilde{\eta}$ . Denote by  $P$  a permutation matrix that  
 exchanges the value of a vector's first and last dimension, then we have the following estimator  $\tilde{\eta}'$ :

$$\tilde{\eta}'(\mathbf{u}) = \begin{cases} \tilde{\eta}(\mathbf{u}), & u_{K+1} \neq \max_y u_y, \\ P\tilde{\eta}(P\mathbf{u}), & \text{else.} \end{cases}$$

466 We then prove that it is a valid estimator. Denote by  $\Delta_+^{K+1}$  the set of all the  $\beta \in \Delta^{K+1}$  that  
 467  $\beta_{K+1} \neq \max_y \beta_y$ , we can learn that  $\Delta_+^{K+1} \in \tilde{\Delta}^{K+1}$ . Then for any  $\beta \in \Delta_+^{K+1}$ , denote by  $\mathbf{u}^*$   
 468 the minimizer of conditional risk w.r.t.  $\beta$  and  $\phi$ , we can learn that  $u_{K+1}^* \neq \max_y u_y^*$  due to the  
 469 **consistency** of  $\phi$  and thus  $\tilde{\eta}'(\mathbf{u}^*) = \tilde{\eta}(\mathbf{u}^*) = \beta$ .

470 For any  $\beta \notin \Delta_+^{K+1}$ , we can learn that  $P\beta \in \Delta_+^{K+1}$ . Denote by  $\mathbf{u}^*$  the minimizer of conditional  
 471 risk w.r.t.  $\beta$  and  $\phi$ , then we can learn that  $P\mathbf{u}^*$  must be the minimizer of  $P\beta$  according to the  
 472 **symmetry** of  $\phi$ . Then we have that  $\tilde{\eta}(P\mathbf{u}^*) = P\beta$ . Furthermore, notice that  $PP\beta = \beta$ , then we  
 473 have  $\tilde{\eta}'(\mathbf{u}^*) = P\tilde{\eta}(P\mathbf{u}^*) = PP\beta = \beta$ .

474 Combining the two paragraphs above, we can learn that  $\tilde{\eta}'$  is a valid multiclass probability estimator  
 475 for  $\phi$ . Then we can construct an unbounded L2D estimator as in (3), which indicates the existence of  
 476 an unbounded probability estimator. Furthermore, since the loss function is unchanged, the collection  
 477 of all the minimizers of  $R_{L_\phi}(\mathbf{u}, \boldsymbol{\eta}(\mathbf{x}), \Pr(M = Y|X = \mathbf{x}))$  are also unchanged, and thus the all the  
 478 values should have the same value on  $\mathcal{U}$ .  $\square$

## 479 B Proof of Proposition 1

*Proof.* The boundedness of the first  $K$  dimensions is straightforward. The  $K + 1$ th dimension's  
 boundedness can also be directly proved by reformulating it as

$$\frac{\exp(u_{K+1})}{\exp(u_{K+1}) + \sum_{y'=1}^K \exp(u_{y'}) - \max_{y' \in \{1, \dots, K\}} \exp(u_{y'})}.$$

Then we begin to prove its maxima-preserving. Denote by  $y_1 = \operatorname{argmax}_{y \in \{1, \dots, K\}} u_y$ . It is easy to  
 verify that  $\tilde{\psi}_{y_1}(\mathbf{u}) > \tilde{\psi}_y(\mathbf{u})$  for any  $y \in \{1, \dots, K\} \setminus \{y_1\}$  due to the property of softmax function.  
 Then we focus on the relation between  $\tilde{\psi}_{y_1}(\mathbf{u})$  and  $\tilde{\psi}_{K+1}(\mathbf{u})$ . We can learn that:

$$\frac{\exp(u_{K+1})}{\exp(u_{K+1}) + \sum_{y'=1}^K \exp(u_{y'}) - \exp(u_{y_1})} = \frac{1}{1 + \frac{\sum_{y'=1}^K \exp(u_{y'}) - \exp(u_{y_1})}{\exp(u_{K+1})}},$$

$$\frac{\exp(u_{y_1})}{\sum_{y'=1}^K \exp(u_{y'})} = \frac{1}{1 + \frac{\sum_{y'=1}^K \exp(u_{y'}) - \exp(u_{y_1})}{\exp(u_{y_1})}}.$$

480 When  $u_{K+1} > u_{y_1}$ , we can learn that  $\frac{\sum_{y'=1}^K \exp(u_{y'}) - \exp(u_{y_1})}{\exp(u_{K+1})} < \frac{\sum_{y'=1}^K \exp(u_{y'}) - \exp(u_{y_1})}{\exp(u_{y_1})}$ , then  
 481  $\tilde{\psi}_{K+1}(\mathbf{u}) > \tilde{\psi}_{y_1}(\mathbf{u})$ . Based on the formulations above, it is also easy to verify the cases when  
 482  $u_{y_1} \geq u_{K+1}$ . Combining the discussions above and we can conclude the proof.  $\square$

483 **C Proof of Theorem 2**

484 *Proof.* The consistency can be directly obtained by recovering the class probability and expert  
 485 accuracy according to the maxima-preserving, then we first focus on how our proposed estimator  
 486 recovers the posterior probabilities.

487 According to the property of log loss, it is easy to learn that  $\tilde{\psi}_y(\mathbf{g}^*(\mathbf{x})) = p(y|\mathbf{x})$  for  $y \in \{1, \dots, K\}$ .  
 488 When  $y = K + 1$ , we can learn from the range of our estimator and the property of binary log loss  
 489 that  $\tilde{\psi}_y(\mathbf{g}^*(\mathbf{x})) = \Pr(M = Y|X = \mathbf{x})$ . Then we can conclude that  $\tilde{\psi}(\mathbf{g}^*(\mathbf{x})) = \boldsymbol{\eta}(\mathbf{x}) \times \Pr(M =$   
 490  $Y|X = \mathbf{x})$ .

491 Then we begin to prove that there is no unbounded probability estimator by contradiction. Suppose  
 492 there exists an unbounded estimator  $\psi$ . For any  $\mathbf{g}$ , it must be the solution of a distribution and expert  
 493 whose posterior probability is  $\tilde{\psi}(\mathbf{g}(\mathbf{x}))$  for each point  $\mathbf{x}$ . However, for a  $\mathbf{g}$  that there exists  $\mathbf{x}$  that  
 494  $\psi(\mathbf{g}(\mathbf{x})) \notin \Delta^K \times [0, 1]$ , we can learn that it cannot be the solution of any distribution and expert  
 495 according to the definition of probability estimator. We can learn from this contradiction that  $\psi$  is not  
 496 a probability estimator as long as its range is not  $\Delta^K \times [0, 1]$ .  $\square$

497 **D Proof of Corollary 1**

*Proof.* We can set the multi-class loss to  $\phi_{\tilde{\psi}}$  to get our proposed loss:

$$\phi_{\tilde{\psi}}(\mathbf{u}, y) = \begin{cases} -\log\left(\frac{\exp(u_y)}{\sum_{y'=1}^{K'-1} \exp(u_{y'})}\right), & y \neq K', \\ -\log\left(\frac{\exp(u_{K'})}{\sum_{y'=1}^{K'} \exp(u_{y'}) - \max_{y' \in \{1, \dots, K'-1\}} \exp(u_{y'})}\right), & \text{else.} \end{cases}$$

A similar result can be deduced for the OvA-based surrogate by considering the following consistent  
 multi-class loss with a strictly proper binary composite loss  $\xi$ :

$$\phi_{\text{OvA}}(\mathbf{u}, y) = \begin{cases} \xi(u_y) + \sum_{y' \neq y, K'} \xi(-u_{y'}), & y \neq K', \\ \xi(u_{K'}) - \xi(-u_{K'}), & \text{else.} \end{cases}$$

498 We then begin to prove their consistency. It is easy to verify that  $\phi_{\tilde{\psi}}$  is minimized when  $\tilde{\psi}(\mathbf{g}(\mathbf{x})) =$   
 499  $\frac{p(y|\mathbf{x})}{1-p(K'|\mathbf{x})}$ , and we can conclude its consistency using the maxima-preserving property of  $\tilde{\psi}$ .

500 For  $\phi_{\text{OvA}}$ , denote by  $\psi_{\text{OvA}}$  the inverse link of  $\xi$ . We can learn that  $\mathbf{g}^*(\mathbf{x}) = \psi_{\text{OvA}}\left(\frac{p(y|\mathbf{x})}{1-p(K'|\mathbf{x})}\right)$ , and  
 501 we can learn the consistency since  $\xi$  is strictly proper and thus  $\psi_{\text{OvA}}$  is increasing.  $\square$

502 **E Proof of Theorem 3**

*Proof.* We first apply Pinsker's inequality. We can learn that:

$$R_{L_{\tilde{\psi}}}(\mathbf{g}) - R_{L_{\tilde{\psi}}}^* \geq \mathbb{E}_{p(\mathbf{x})} \left[ \frac{1}{2} \|\tilde{\psi}_{1:K}(\mathbf{g}(\mathbf{x})) - \boldsymbol{\eta}(\mathbf{x})\|_1^2 + 2(\tilde{\psi}_{K+1}(\mathbf{g}(\mathbf{x})) - \Pr(M = Y|X = \mathbf{x}))^2 \right].$$

503 We can learn that  $R_{L_{\tilde{\psi}}}(\mathbf{g}) - R_{L_{\tilde{\psi}}}^* \geq \mathbb{E}_{p(\mathbf{x})} \left[ \frac{1}{2} \|\tilde{\psi}_{1:K}(\mathbf{g}(\mathbf{x})) - \boldsymbol{\eta}(\mathbf{x})\|_1^2 \right]$  immediately and learn that  
 504 the bound for  $R_{01}$  following the analysis of cross-entropy loss in ordinary classification. Then we  
 505 move to analyze the 0-1-deferral risk. We can further learn that :

$$\begin{aligned} R_{L_{\tilde{\psi}}}(\mathbf{g}) - R_{L_{\tilde{\psi}}}^* &\geq \mathbb{E}_{p(\mathbf{x})} \left[ \frac{1}{2} \|\tilde{\psi}_{1:K}(\mathbf{g}(\mathbf{x})) - \boldsymbol{\eta}(\mathbf{x})\|_1^2 + 2(\tilde{\psi}_{K+1}(\mathbf{g}(\mathbf{x})) - \Pr(M = Y|X = \mathbf{x}))^2 \right] \\ &\geq \frac{1}{2} \mathbb{E}_{p(\mathbf{x})} \left[ \|\tilde{\psi}_{1:K}(\mathbf{g}(\mathbf{x})) - \boldsymbol{\eta}(\mathbf{x})\|_1^2 + (\tilde{\psi}_{K+1}(\mathbf{g}(\mathbf{x})) - \Pr(M = Y|X = \mathbf{x}))^2 \right] \end{aligned}$$

506 For any  $\mathbf{x}$ , when  $\mathbf{g}(\mathbf{x})$  can induce the Bayes optimal solution for it, the excess risk is zero and the  
 507 bound holds naturally. When it is not optimal, denote by  $\eta_{K+1}(\mathbf{x}) = \Pr(M = Y|X = \mathbf{x})$ .  $y'$  is the

508 dimension with the largest value of  $\mathbf{g}(\mathbf{x})$  and that of  $\boldsymbol{\eta}(\mathbf{x})$  is  $y''$ . Then we can learn that:

$$\begin{aligned} & \frac{1}{2} \left( \|\tilde{\psi}_{1:K}(\mathbf{g}(\mathbf{x})) - \boldsymbol{\eta}(\mathbf{x})\|_1^2 + (\tilde{\psi}_{K+1}(\mathbf{g}(\mathbf{x})) - \Pr(M = Y|X = \mathbf{x}))^2 \right) \\ & \geq \frac{1}{2} \left( \tilde{\psi}_{y'}(\mathbf{g}(\mathbf{x})) - \eta_{y'}(\mathbf{x}) - \tilde{\psi}_{y''}(\mathbf{g}(\mathbf{x})) + \eta_{y''}(\mathbf{x}) \right)^2 \\ & \geq \frac{1}{2} (\eta_{y'}(\mathbf{x}) - \eta_{y''}(\mathbf{x}))^2 \end{aligned}$$

The last step is obtained according to the maxima-preserving property. Further generalizing  $y'$  and  $y''$  to be instance-dependent ( $y'(\mathbf{x})$  and  $y''(\mathbf{x})$ ), we can learn the following inequality using Jensen's inequality:

$$R_{L_{\tilde{\psi}}}(\mathbf{g}) - R_{L_{\tilde{\psi}}}^* \geq \frac{1}{2} (\mathbb{E}_{p(\mathbf{x})} [|\eta_{y'(\mathbf{x})} - \eta_{y''(\mathbf{x})}|])^2,$$

509 which concludes the proof since the second expectation term is  $R_{01}^\perp(\mathbf{g}) - R_{01}^{\perp*}$  □

## 510 F Details of Experiments

511 **Details of Model and Optimizer:** For all the methods on different datasets, we use the 28-layer  
 512 WideResNet that is the same as those used in Mozannar and Sontag [28], Charusaie et al. [9]. The  
 513 optimizer is SGD with cosine annealing, where the learning rate is 1e-1 and weight decay is 5e-4.  
 514 We conduct the experiments on 8 NVIDIA GeForce 3090 GPUs and the batch size is 1024 (128 on  
 515 each GPU). The training epoch on CIFAR100 is set to 200 and 400 on CIFAR10H, respectively.

**Details of Evaluation Metrics:** The reported Error is the sample mean of  $\ell_{01}^\perp$ , and Coverage is the ratio of undeferred samples. The ECE of expert accuracy is defined below:

$$\text{ECE} = \sum_{i=1}^N b_i |p_i - c_i|,$$

516 where  $b_i$  is the ratio of predictions whose confidences fall into the  $i$ th bin.  $p_i$  is the average confidence  
 517 and  $c_i$  is the average accuracy in this bin. We set the bin number to 15. The budgeted error is  
 518 calculated as below: if the coverage is lower than  $1 - x\%$ , we will use the classifier's prediction  
 519 instead of the expert's for those samples whose estimated expert accuracy is lower to make the  
 520 coverage equal to  $1 - x\%$ .

## 521 G Limitations and Broader Impact

522 **Limitations:** This work is designed for L2D without extra constraints on the number of expert  
 523 queries. We believe that combining it with selective learning, i.e., adding explicit constraints on the  
 524 ratio of deferred samples, can be a promising future direction.

525 **Broader Impact:** When applied in real-world applications, the frequency of expert queries may  
 526 be imbalanced due to the performance differences of the expert among samples. This is a common  
 527 impact shared by all the L2D methods. We believe that introducing fairness targets into L2D can be  
 528 another promising direction.