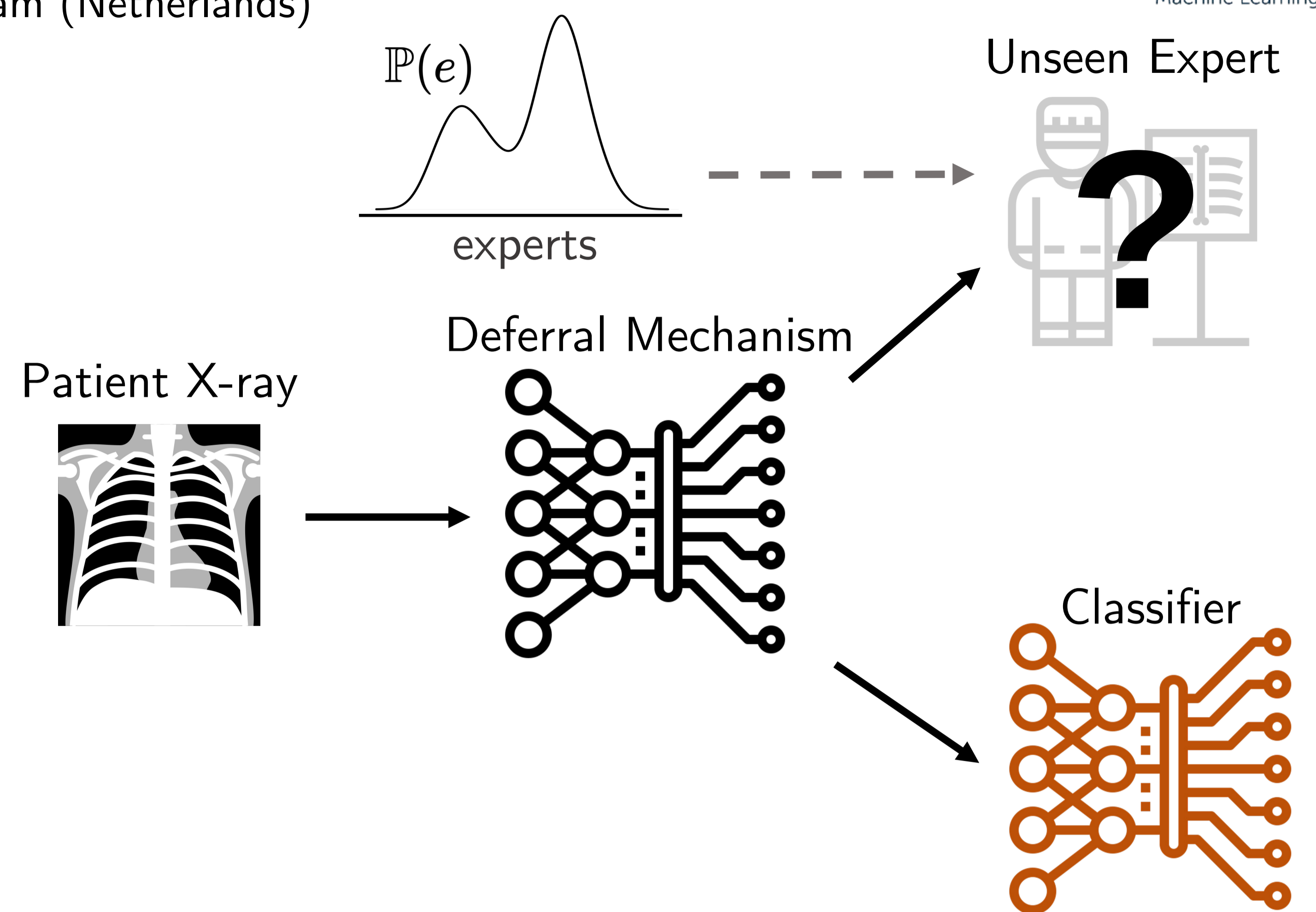


Learning-to-defer (L2D) is a proposal for **hybrid intelligent** systems that gives the AI the option to abstain and defer its prediction to a human upon facing a challenging or high-risk decision. Existing L2D systems are trained to be customized to one (or more) specific humans and if the expert were to change, the system should be re-trained.

We formulate an L2D system that can cope with **unseen experts** at test-time by training its deferral subcomponent to generalize to all experts in a population. We propose a general **meta-learning** implementation that can adapt to any expert by only using a **context set** of demonstrations.



## Background: L2D [Mozannar & Sontag, 2020]

The learning problem involves jointly training two sub-models: a classifier  $h: \mathcal{X} \rightarrow \mathcal{Y}$  and a rejector  $r: \mathcal{X} \rightarrow \{0, 1\}$ . When  $r(x) = 0$ , the classifier makes the decision, and when  $r(x) = 1$  the classifier abstains and defers the decision to the human.

$$\mathcal{D} = \{\mathbf{x}_n, \mathbf{y}_n, m_n\}_{n=1}^N$$

Classifier-rejector (0-1) loss:

$$L_{0-1}(h, r) = \mathbb{E}_{\mathbf{x}, \mathbf{y}, m} \left[ \underbrace{(1 - r(\mathbf{x})) \mathbb{I}[h(\mathbf{x}) \neq \mathbf{y}]}_{\text{predict}} + \underbrace{r(\mathbf{x}) \mathbb{I}[m \neq \mathbf{y}]}_{\text{defer}} \right]$$

They propose a reduction from multiclass expert deferral to cost sensitive learning by unifying the classifier and rejector via an augmented label space that includes the rejection option:  $\mathcal{Y}^\perp = \mathcal{Y} \cup \{\perp\}$

Then a (consistent) surrogate loss is constructed that extends cross entropy loss. [Verma & Nalisnick, 2022] also showed one-vs-all parameterization is also a consistent surrogate.

$$\phi_{SM}(g_1, \dots, g_K, g_\perp; \mathbf{x}, \mathbf{y}, m) = -\log \left( \frac{\exp\{g_y(\mathbf{x})\}}{\sum_{y' \in \mathcal{Y}^\perp} \exp\{g_{y'}(\mathbf{x})\}} \right)$$

Rejection function:

$$\mathbb{I}[g_\perp(\mathbf{x}) \geq \max_k g_k(\mathbf{x})] - \mathbb{I}[m = y] \log \left( \frac{\exp\{g_\perp(\mathbf{x})\}}{\sum_{y' \in \mathcal{Y}^\perp} \exp\{g_{y'}(\mathbf{x})\}} \right)$$

## Learning to Defer to a Population "L2D-Pop"

We assume a generative process for experts from which experts can be sampled indefinitely and without repetition.

$$\mathcal{E} \sim \mathbb{P}(\mathcal{E}), \quad m \sim \mathbb{P}(m|\mathbf{x}, \mathbf{y}, \mathcal{E})$$

Formulation resembles single-expert L2D but now rejector also takes as input some representation of the currently-available expert.

$$r: \mathcal{X} \times \mathcal{E} \rightarrow \{0, 1\}$$

Surrogate loss:

$$\phi_{SM-Pop}(g_1, \dots, g_K, g_\perp; \mathbf{x}, \mathbf{y}, \{m_e, \psi_e^\mathcal{E}\}_{e=1}^E) = \sum_{e=1}^E -\log \left( \frac{\exp\{g_y(\mathbf{x})\}}{\mathcal{Z}(\mathbf{x}, \psi_e^\mathcal{E})} \right) - \mathbb{I}[m_e = y] \log \left( \frac{\exp\{g_\perp(\mathbf{x}, \psi_e^\mathcal{E})\}}{\mathcal{Z}(\mathbf{x}, \psi_e^\mathcal{E})} \right)$$

## Applying Single-Expert L2D to L2D-Pop:

We can apply single-expert L2D to the population setting to learn a rejector that models the population's marginal probability of correctness. This requires a reformulation of L2D-Pop surrogate loss where rejector no longer a function of the expert and the sum over experts 'pushes through' to the indicator term.

## References

Mozannar and Sontag. *Consistent Estimators for Learning to Defer to an Expert*. ICML, 2020.

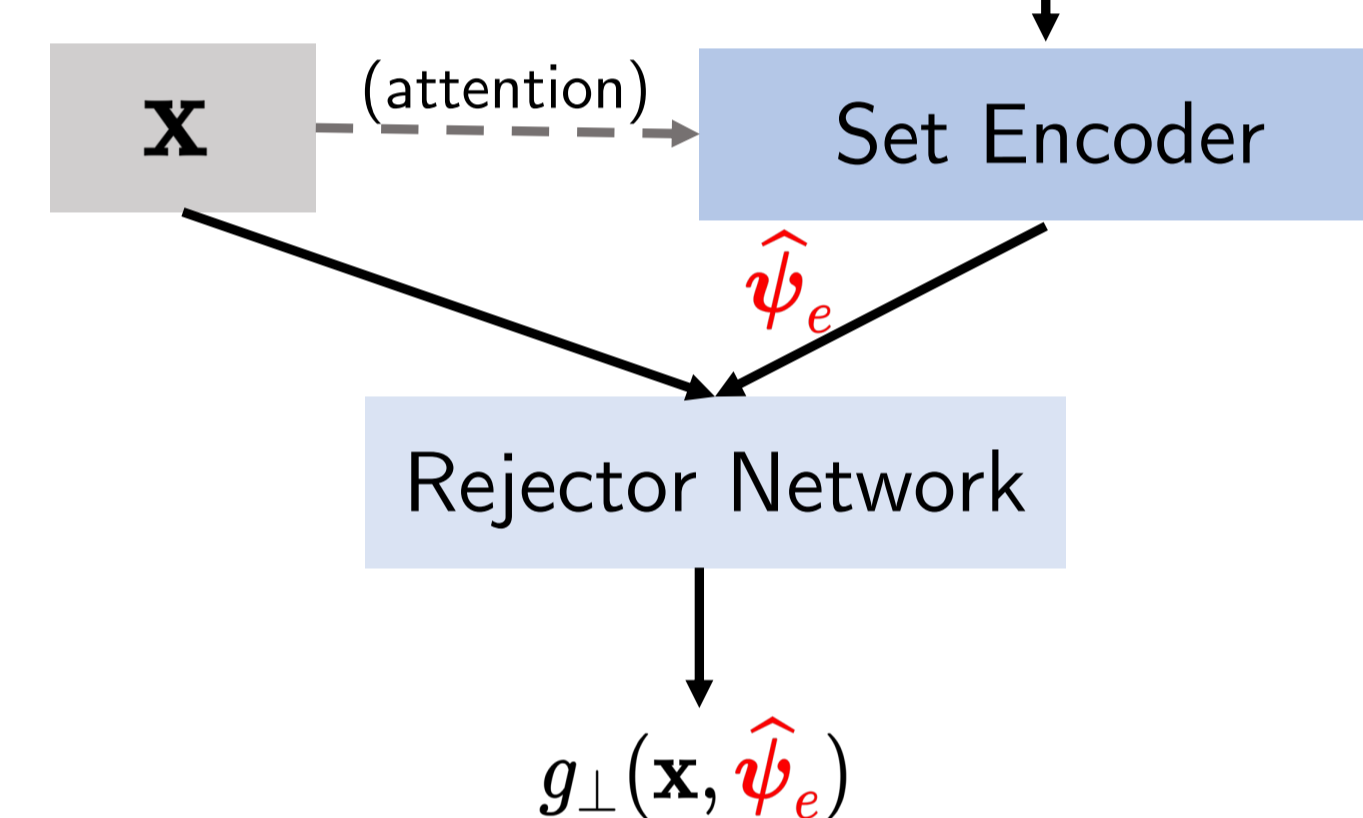
Verma and Nalisnick. *Calibrated Learning to Defer with One-vs-All Classifiers*. ICML, 2022.

## Meta-Learning to Defer

### Model-based approach

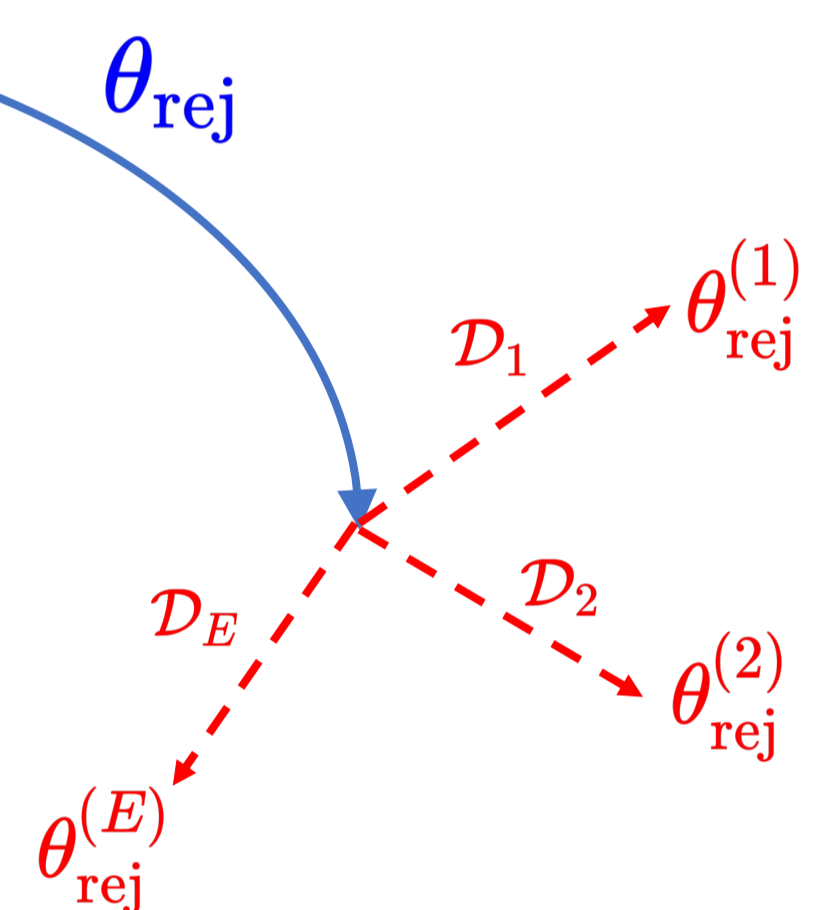
context set / representative set of expert demonstrations

$$\mathcal{D}_e = \{(\mathbf{x}_{e,b}, \mathbf{y}_{e,b}, m_{e,b})\}_{b=1}^B$$



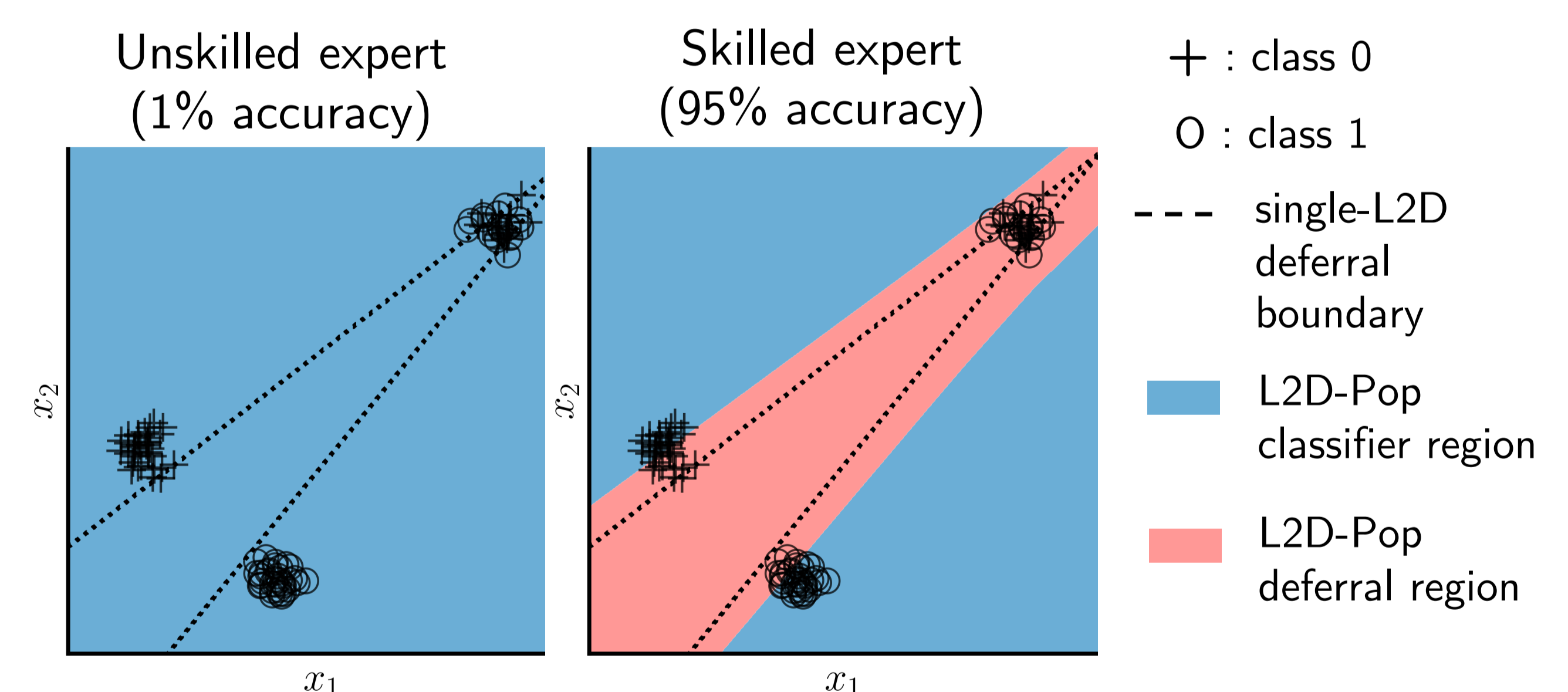
### Meta-optimization approach

- 1 Model the *marginal* expert by the single-expert formulation
- 2 Then finetune on expert context set at test-time



## Experiments

### Synthetic Data



### Varying Population Diversity

