

---

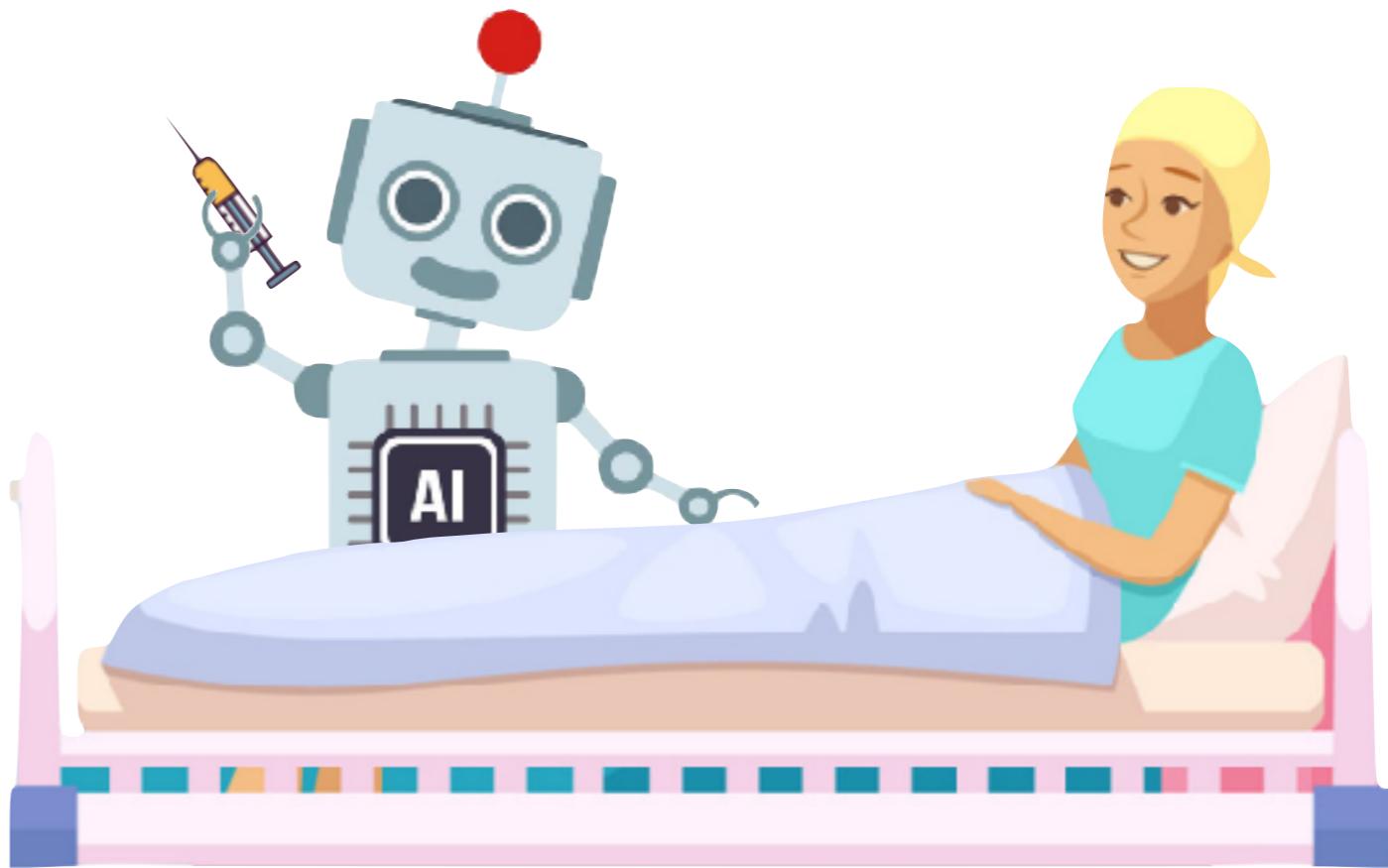
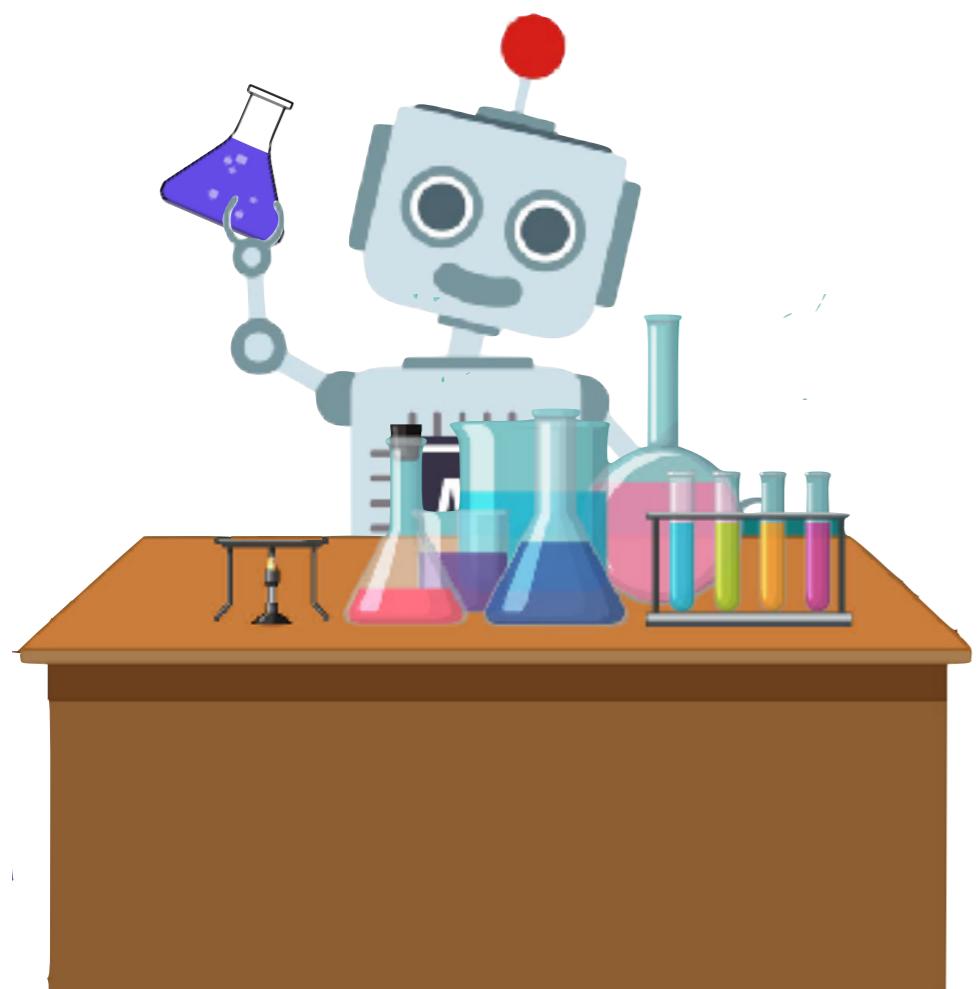
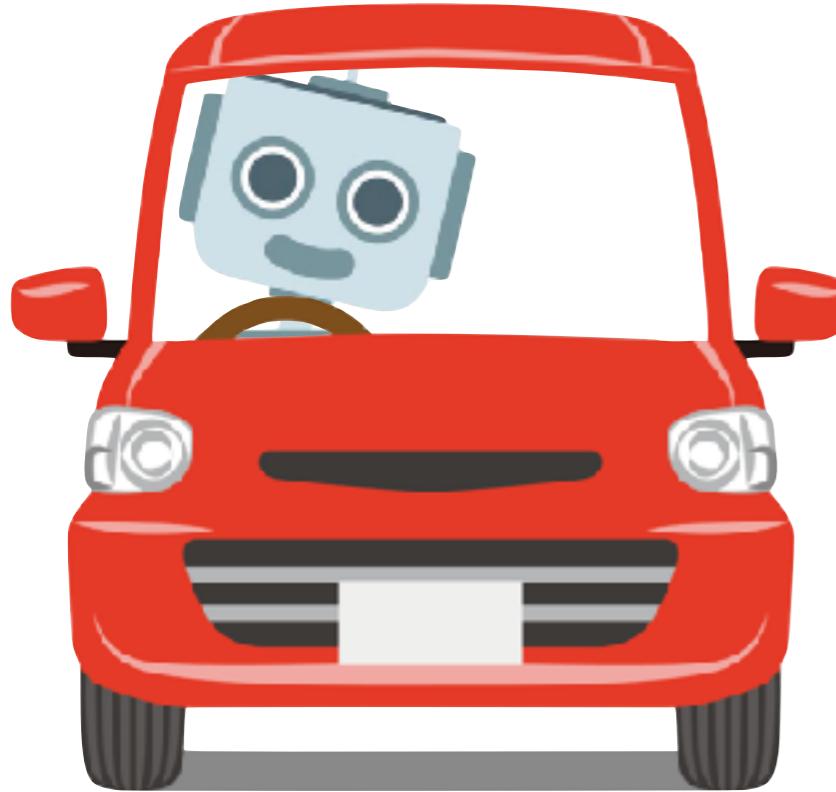
# On the Calibration of Learning-to-Defer Systems

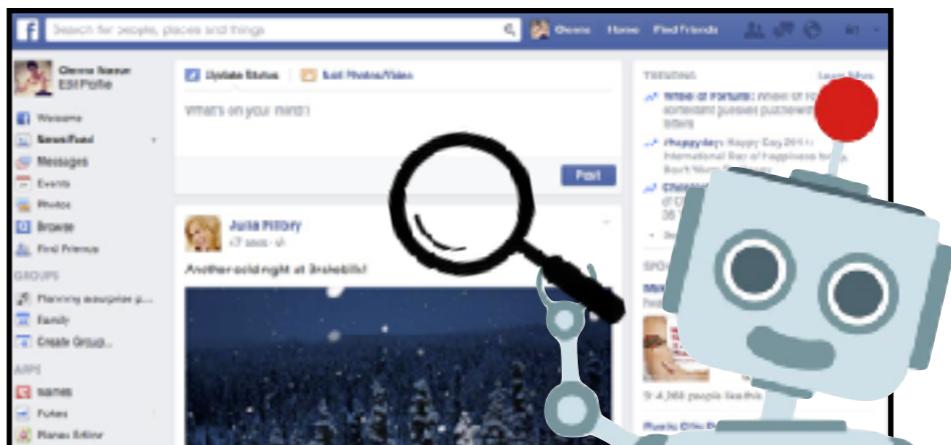
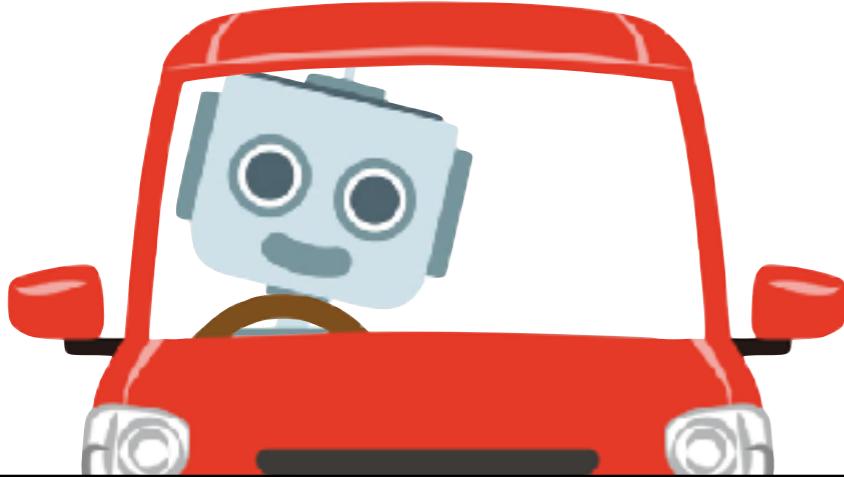
---

Eric Nalisnick



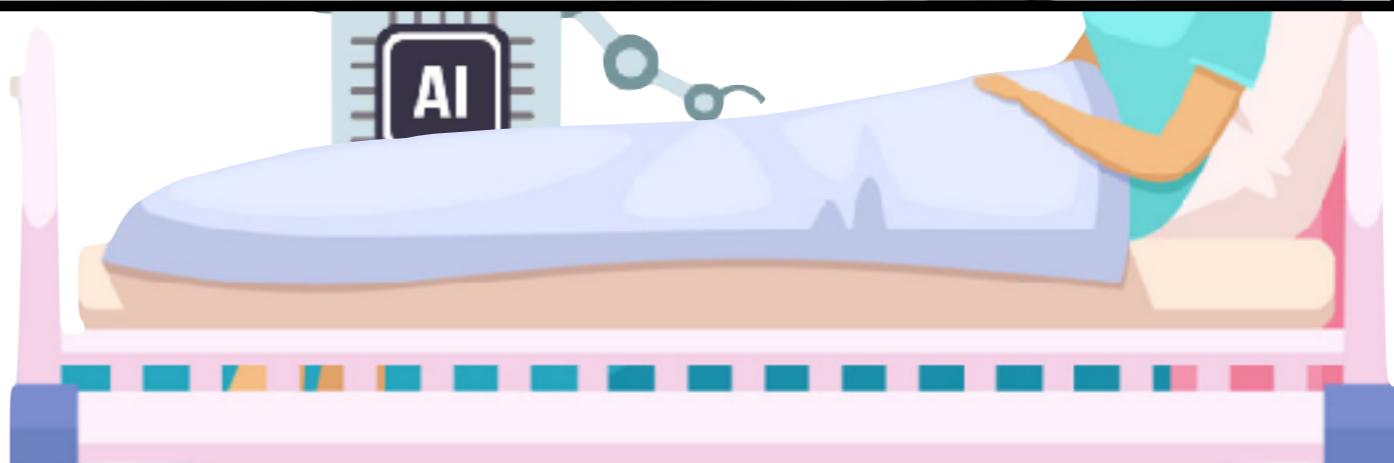
Rajeev Verma

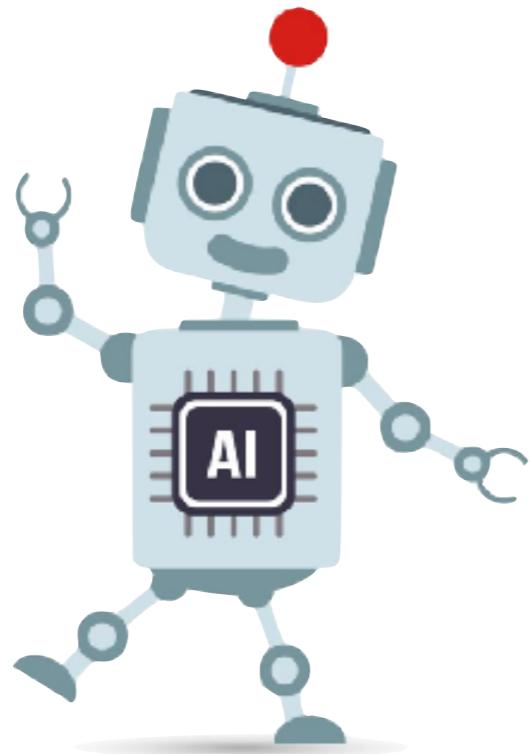




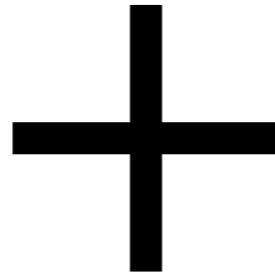
# Tesla recall: 'Full Self-Driving' software runs stop signs

By TOM KRISHER an hour ago





Model



Human

Artificial Intelligence

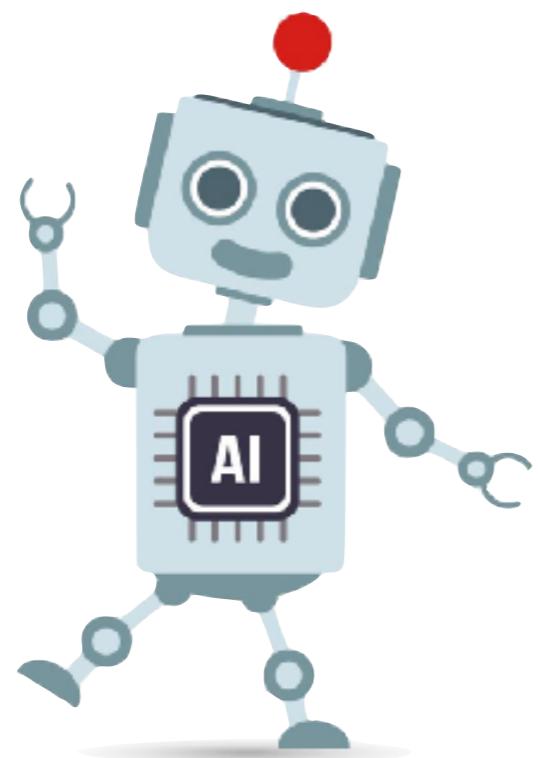
Hybrid Intelligence

## **Junior and senior radiologists benefit from AI assistance when identifying pulmonary nodules**

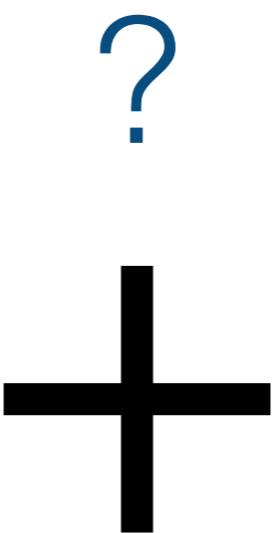
*Hannah Murphy | January 03, 2022 | AI*



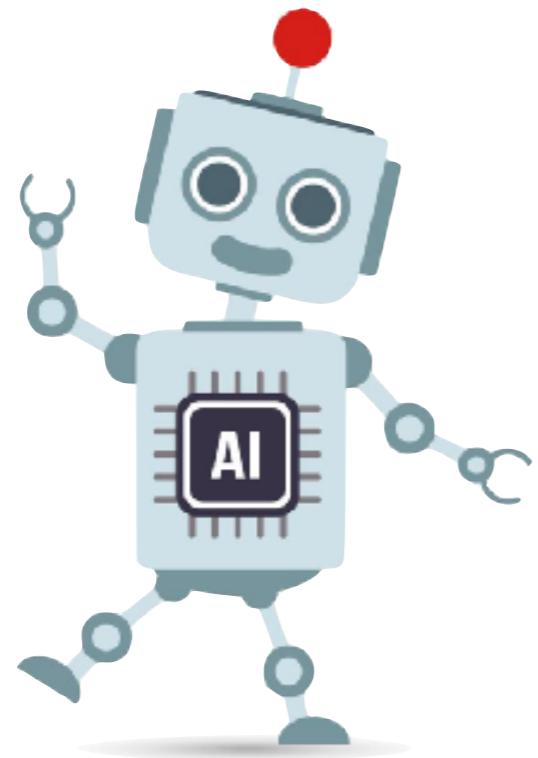
Radiologists benefit from AI-assistance when identifying pulmonary nodules on chest radiographs regardless of their experience level, according to recent research. Since most [early lung cancers](#) initially present as pulmonary nodules, these new findings could benefit patients.



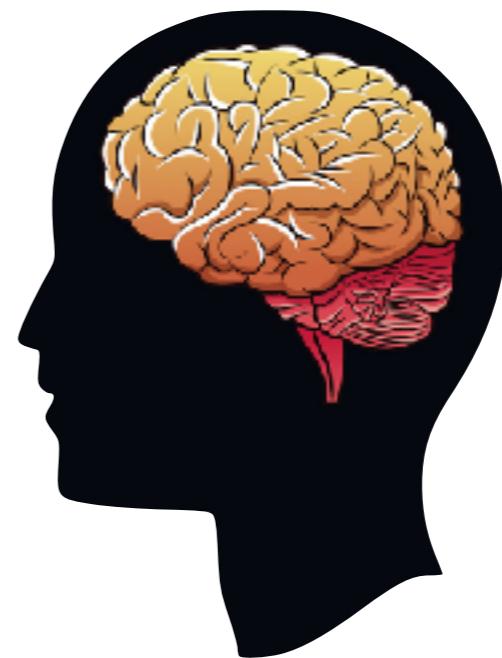
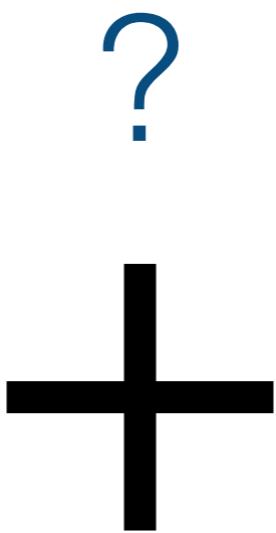
Model



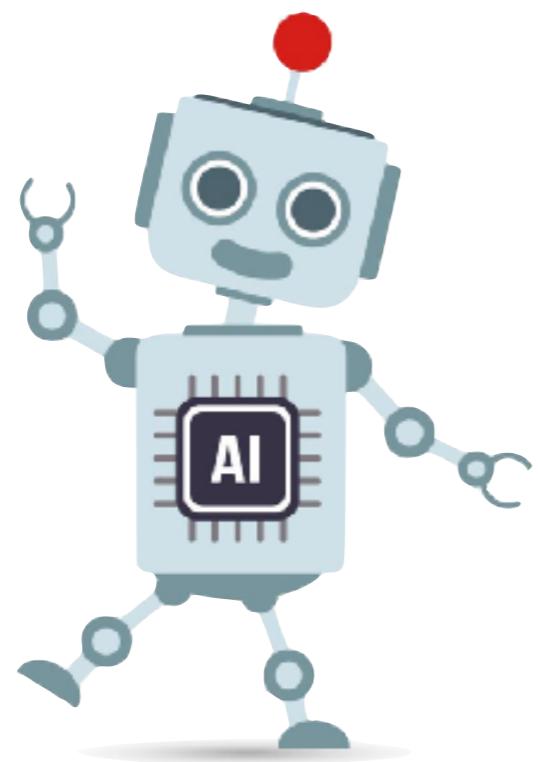
Human



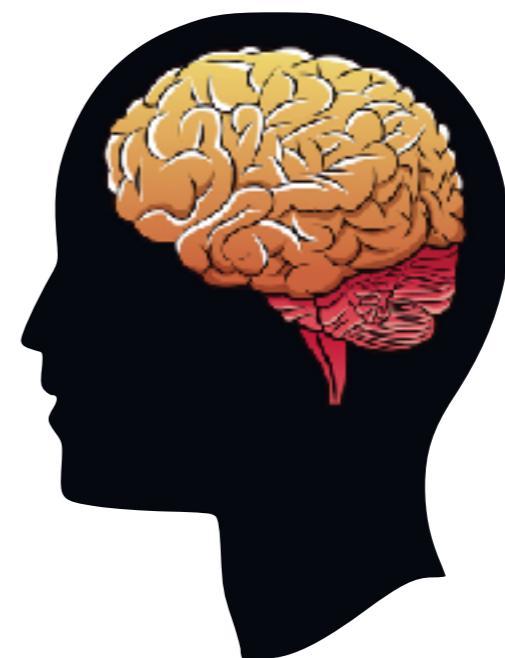
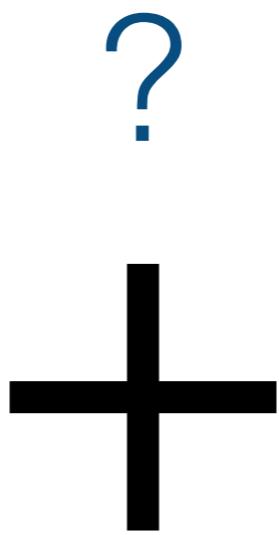
Classifier



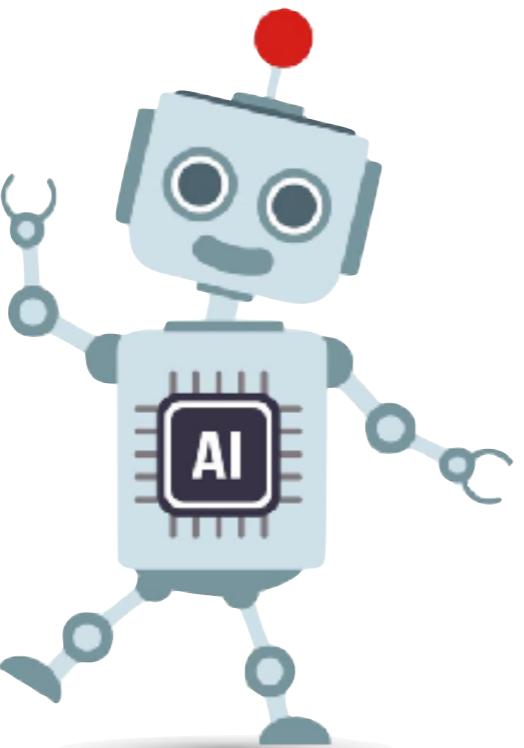
Human



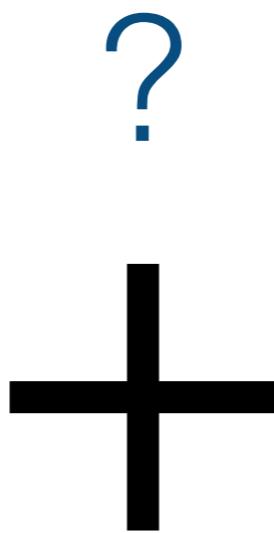
Classifier



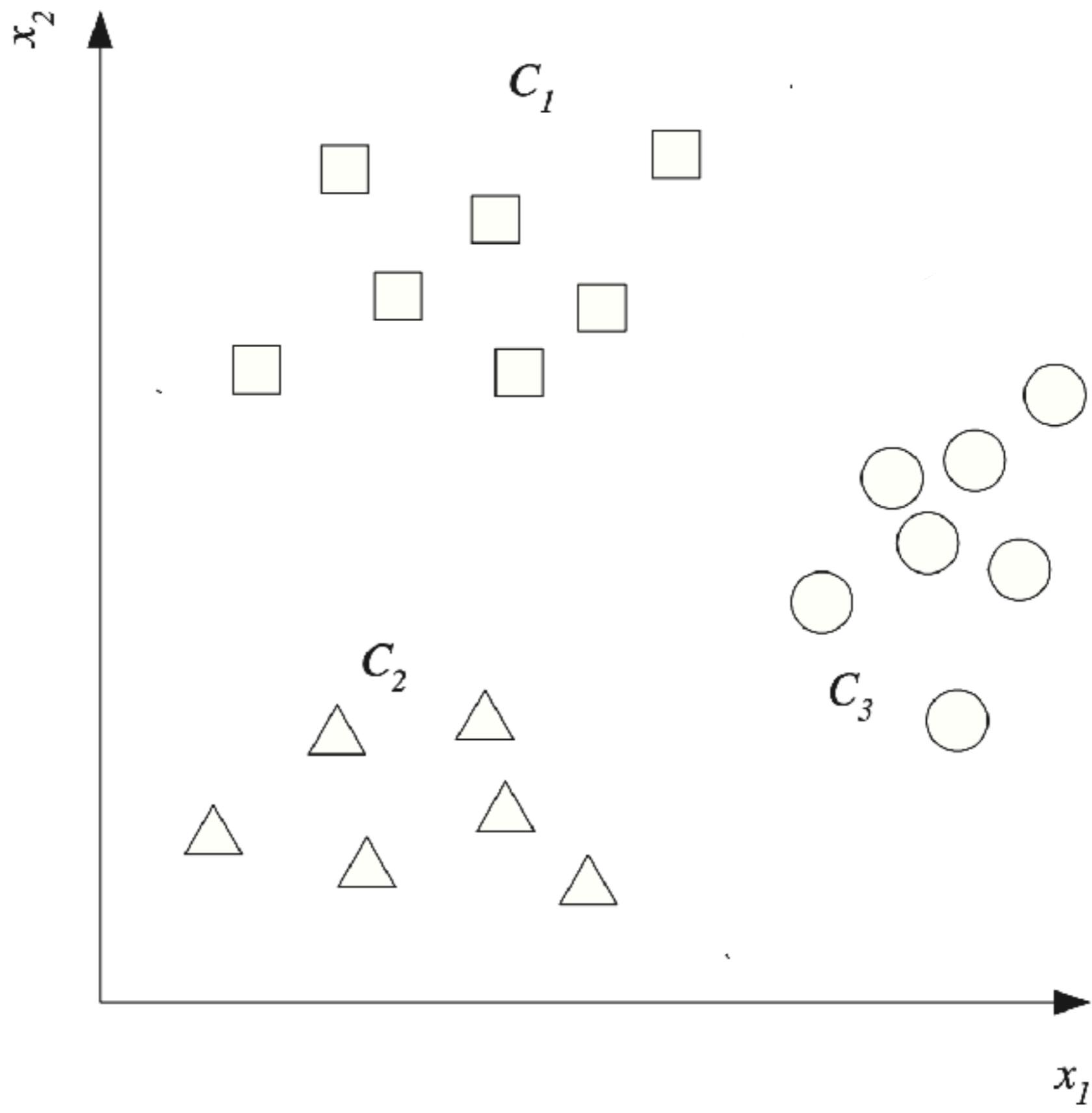
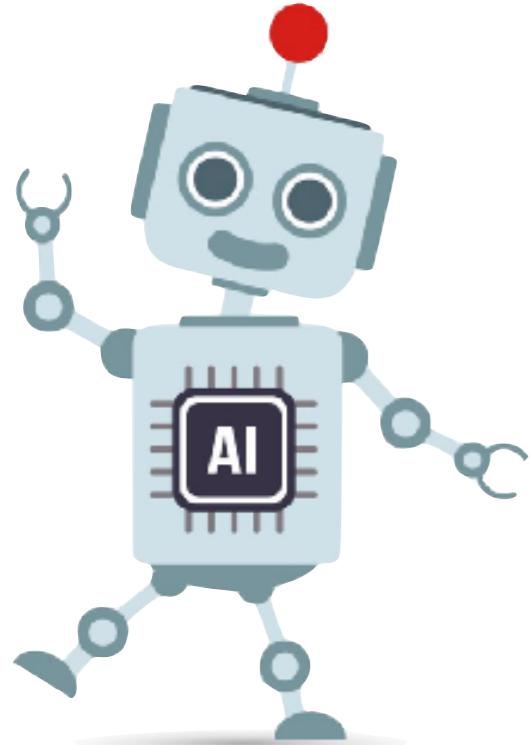
Expert



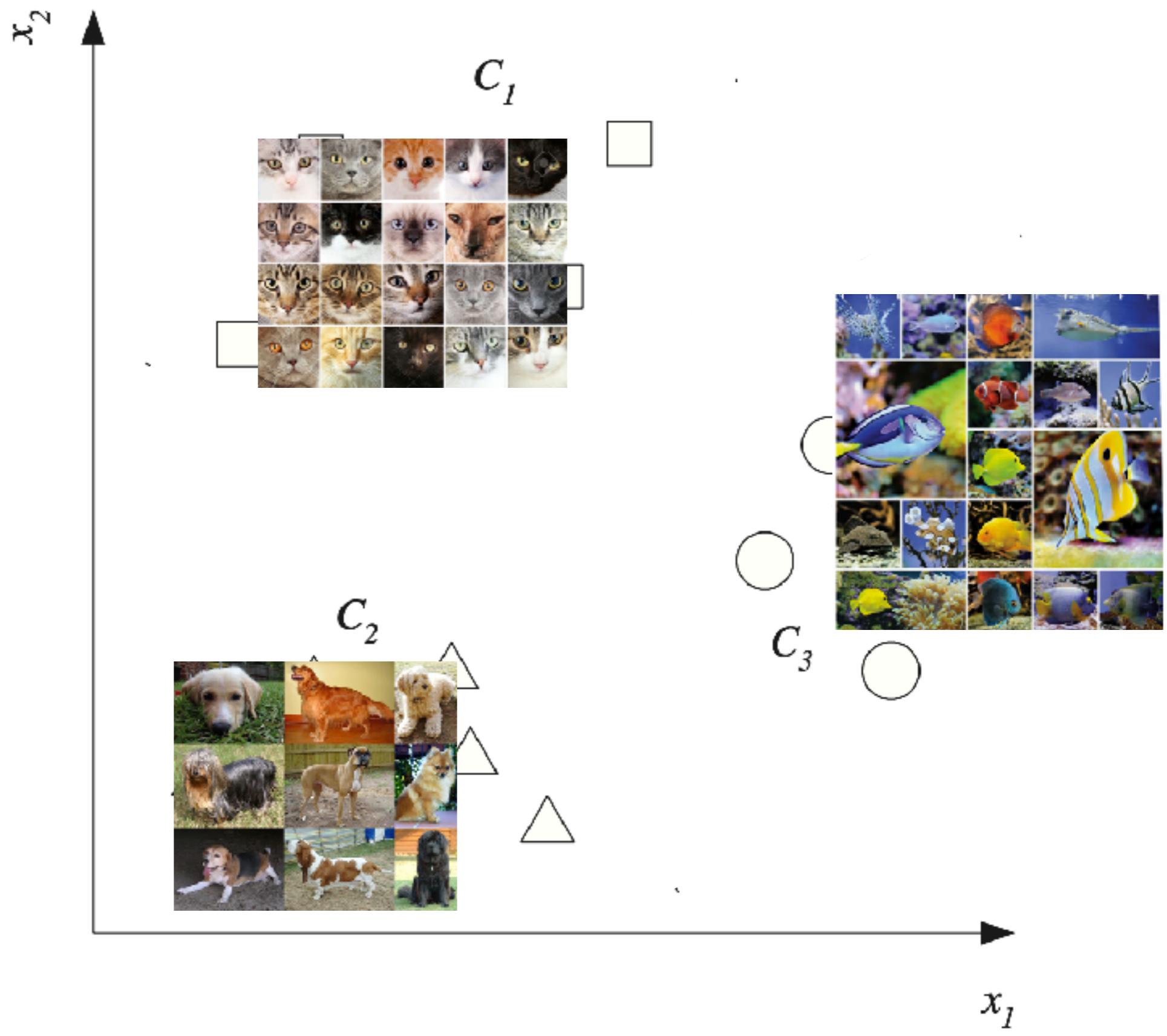
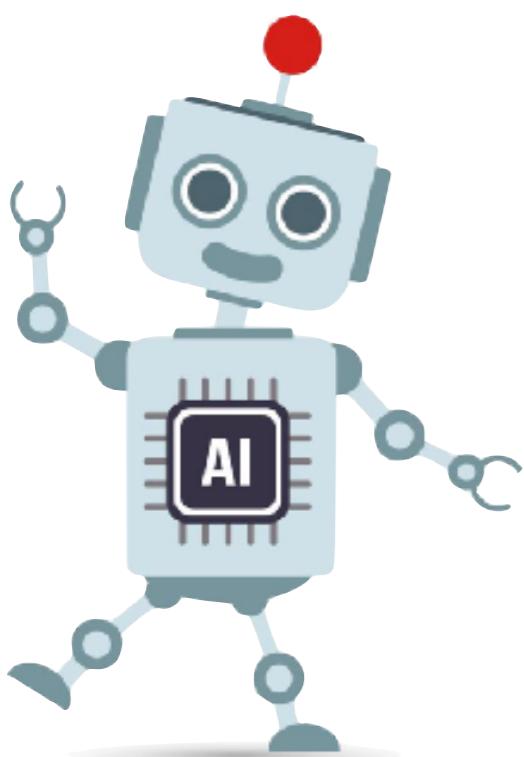
Classifier



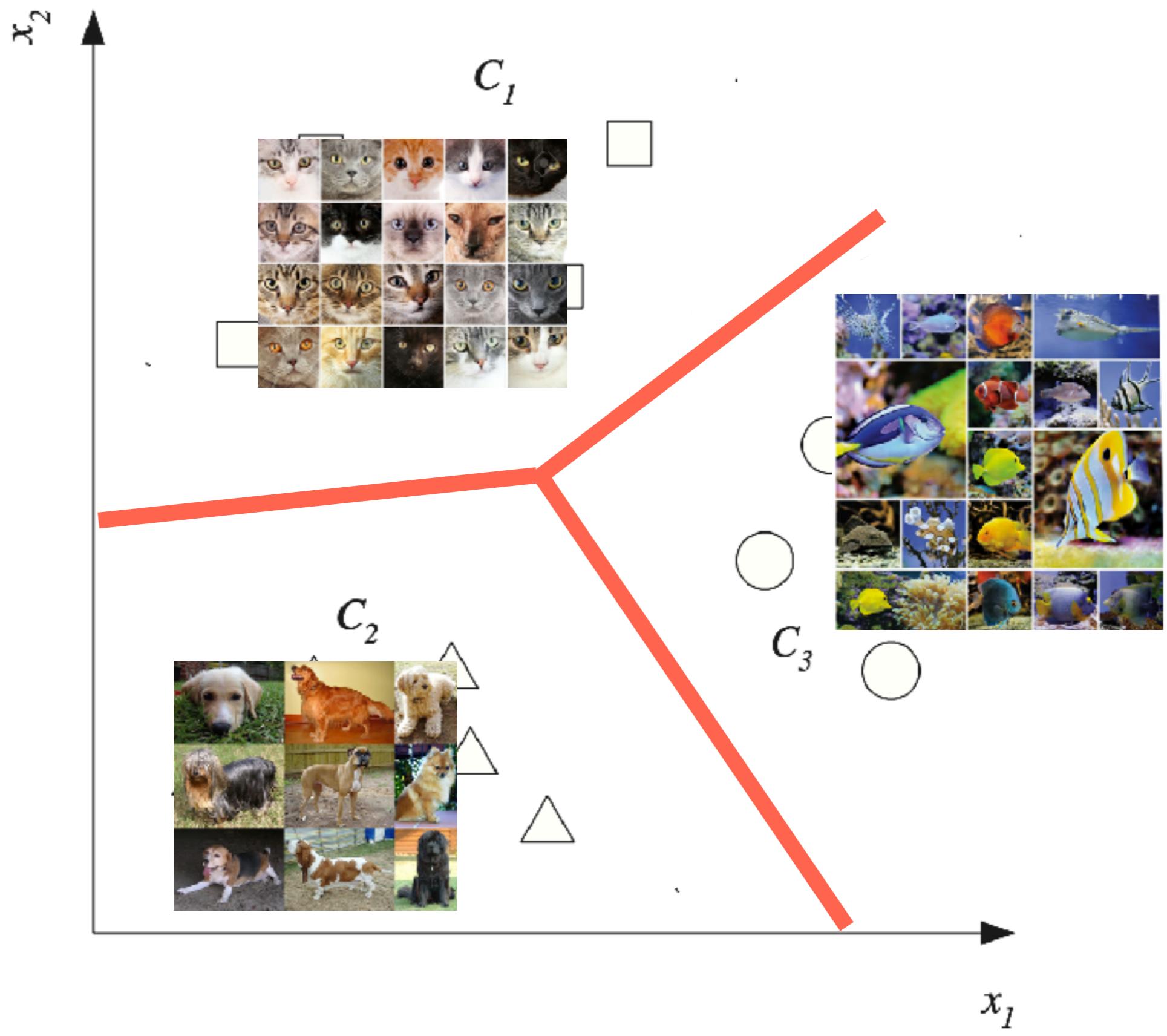
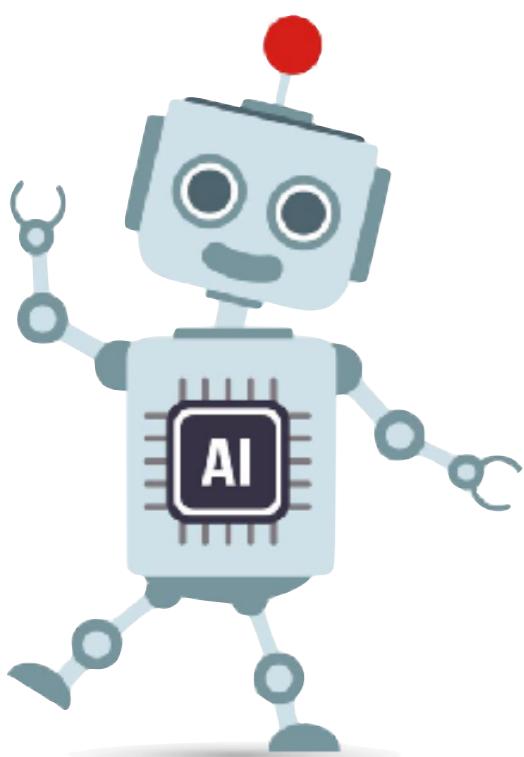
Expert



Classifier

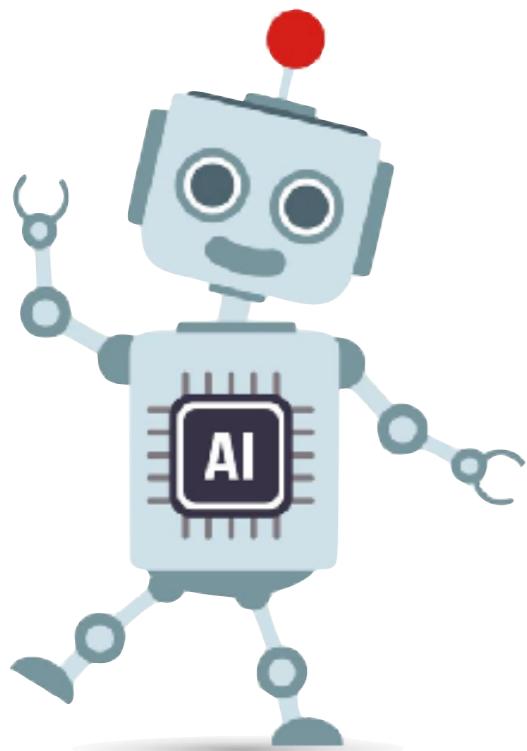


Classifier



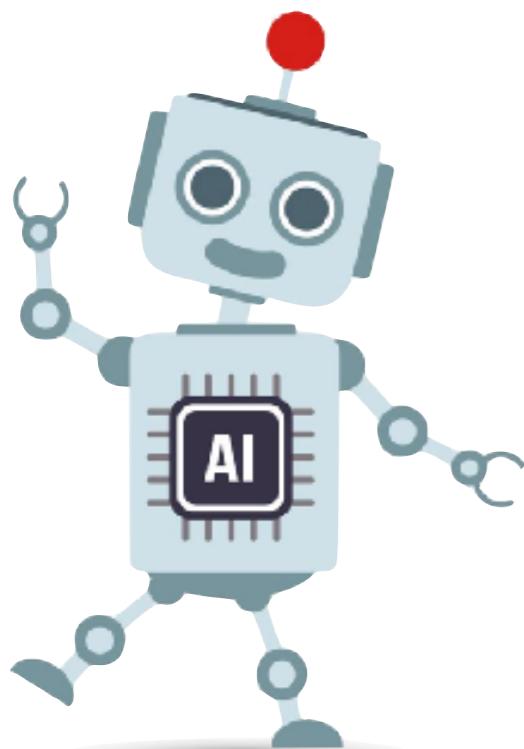
Classifier

**Goal:** For input features  $\mathbf{x}$ ,  
predict membership in one  
of  $K$  classes:  $C_1, \dots, C_K$



Classifier

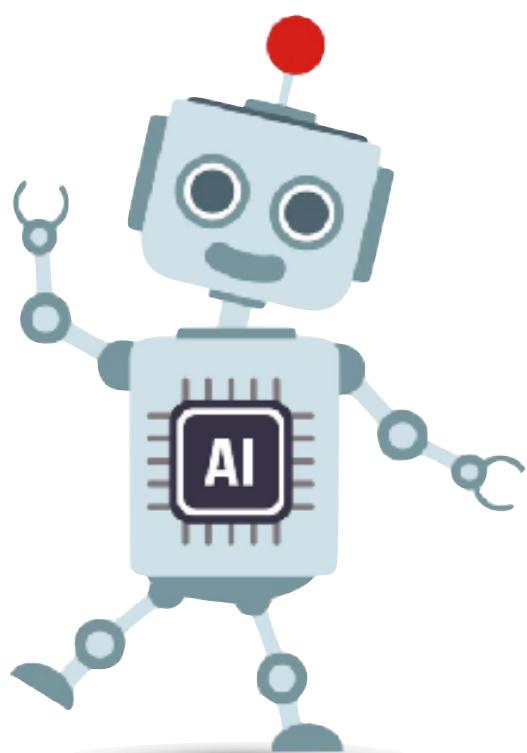
**Goal:** For input features  $\mathbf{x}$ ,  
predict membership in one  
of  $K$  classes:  $C_1, \dots, C_K$



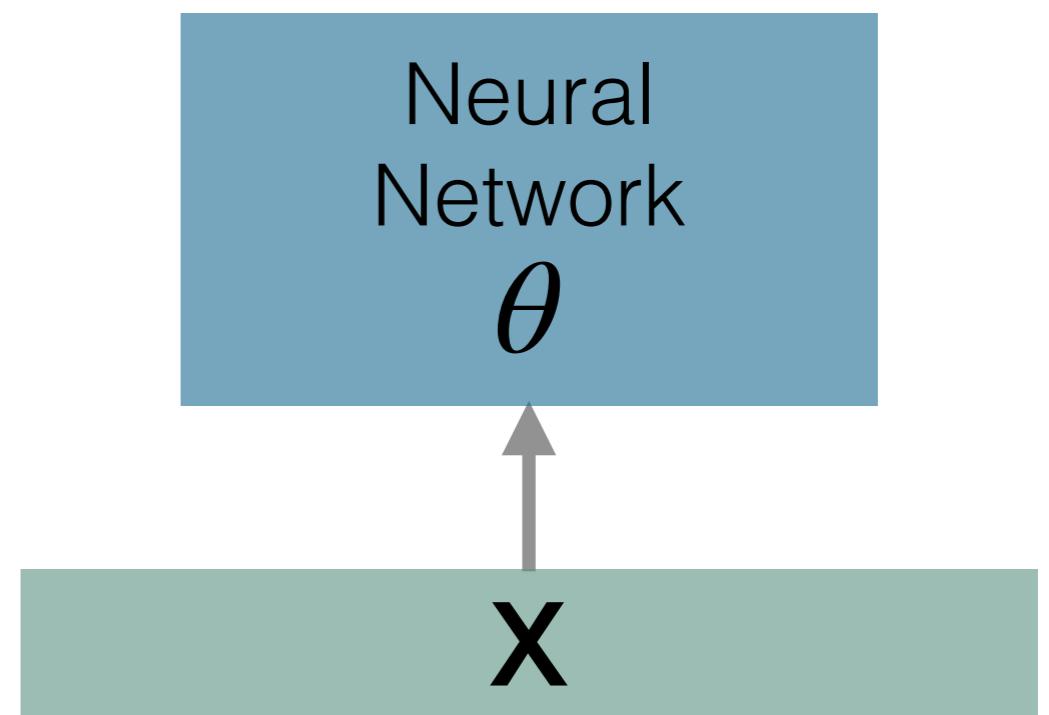
Classifier

$\mathbf{x}$

**Goal:** For input features  $\mathbf{x}$ ,  
predict membership in one  
of  $K$  classes:  $C_1, \dots, C_K$

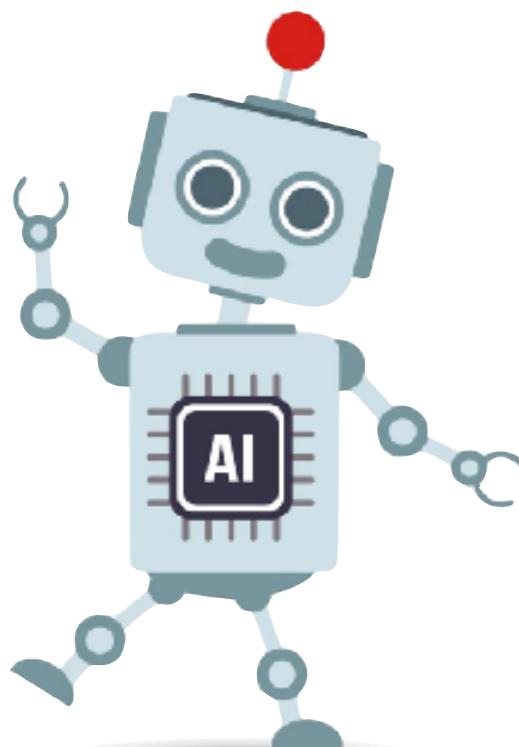


Classifier

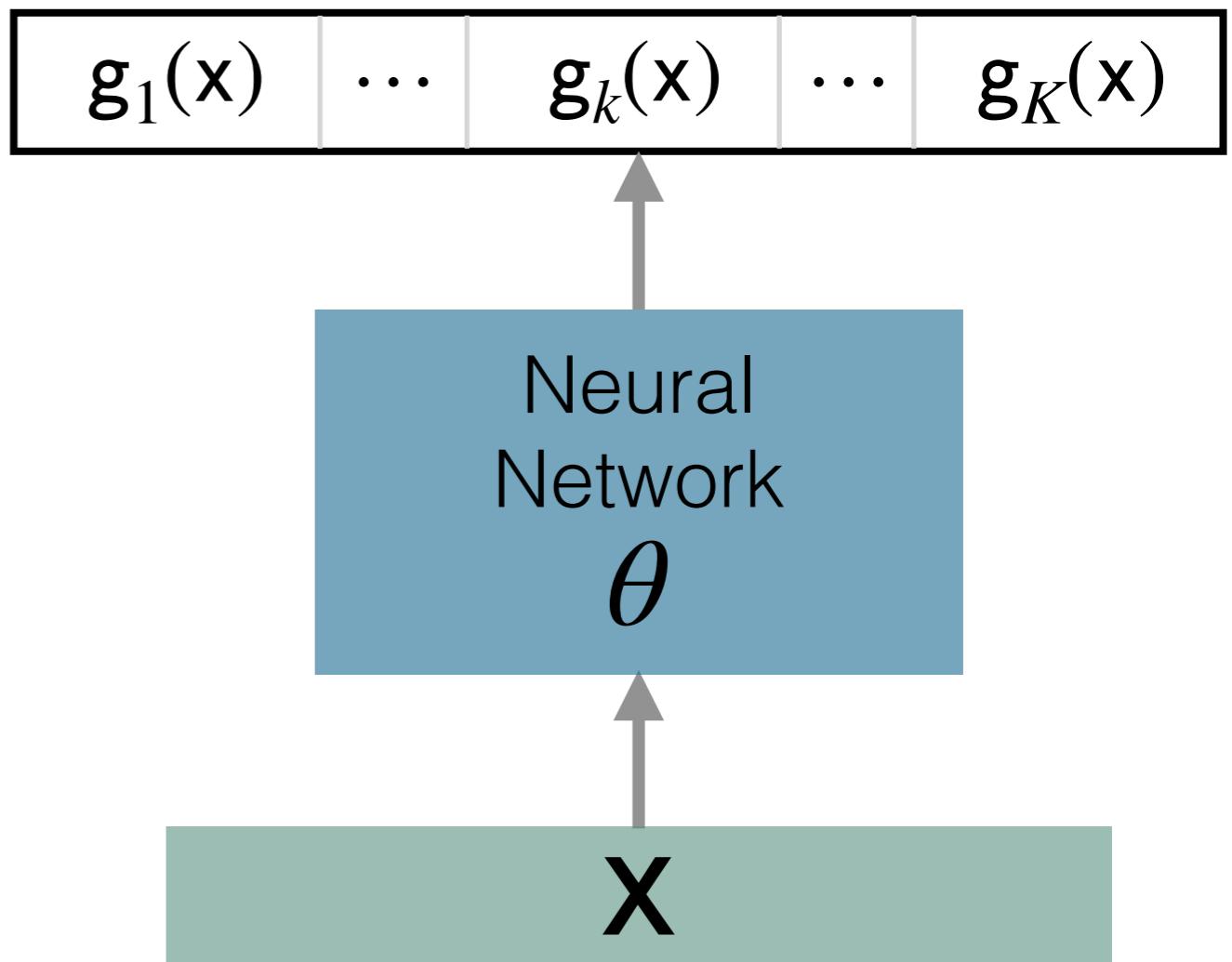


Goal: For input features  $\mathbf{x}$ ,  
predict membership in one  
of  $K$  classes:  $C_1, \dots, C_K$

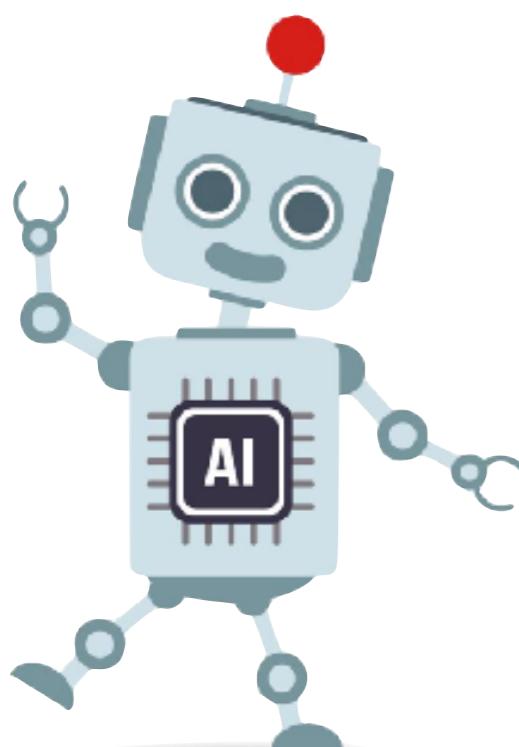
$$g_k(\mathbf{x}) \in \mathbb{R}$$



Classifier

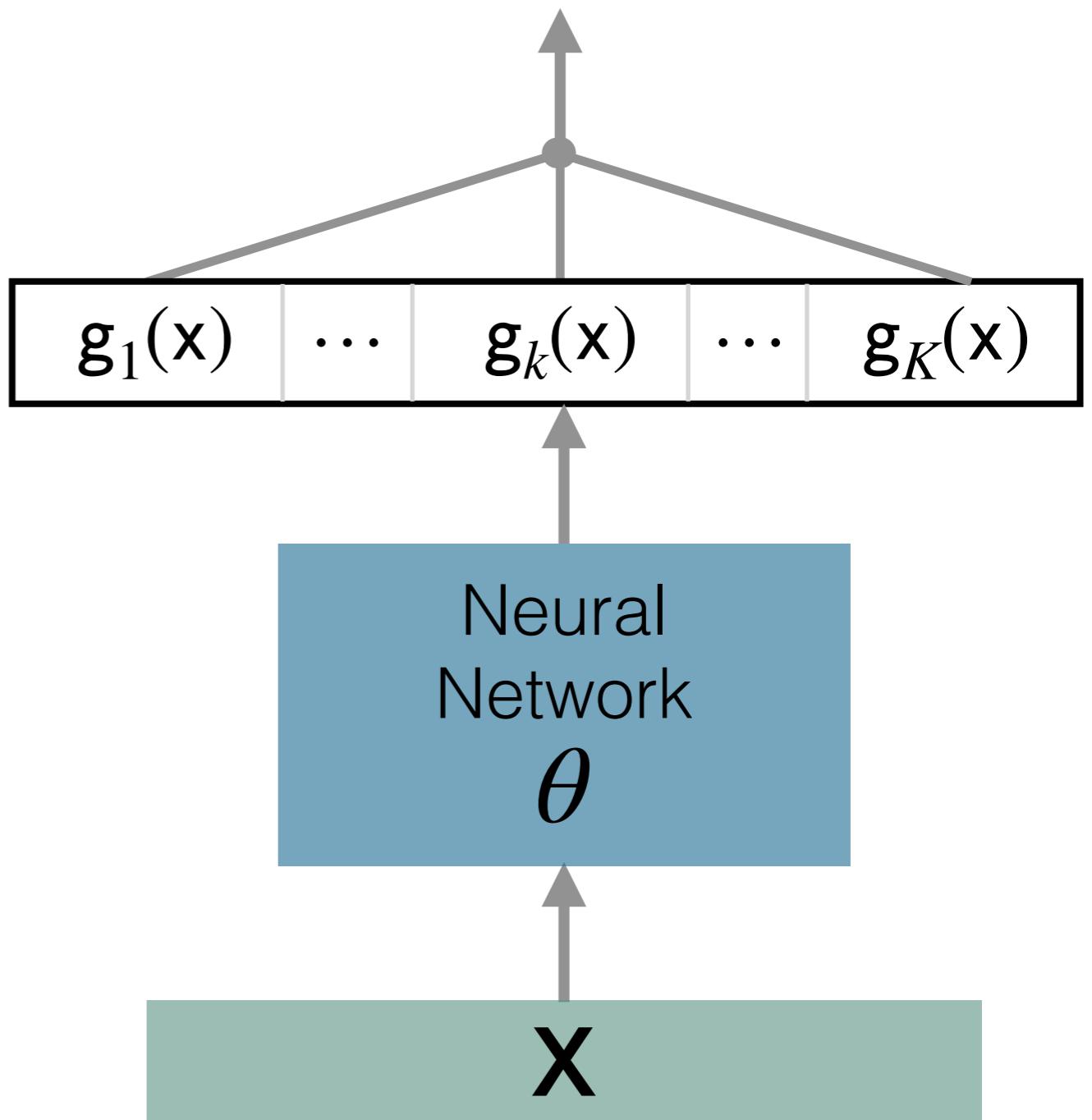


Goal: For input features  $\mathbf{x}$ ,  
predict membership in one  
of  $K$  classes:  $C_1, \dots, C_K$



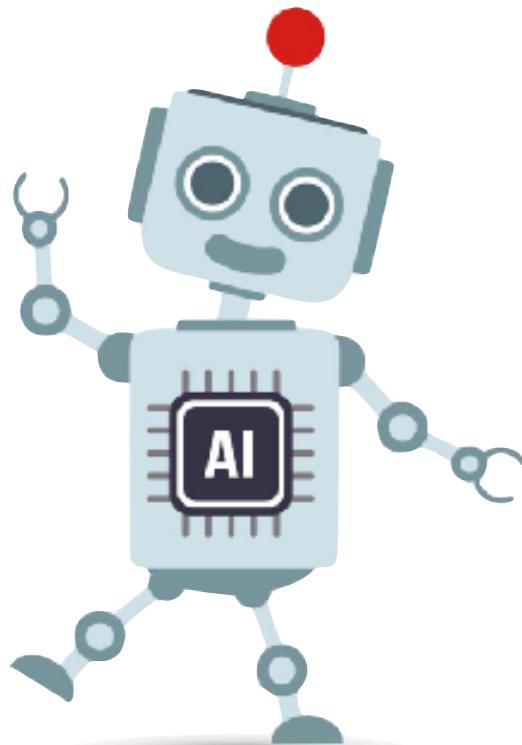
Classifier

$$P(C_i | \mathbf{x}) = \frac{\exp\{g_i(\mathbf{x})\}}{\sum_{k=1}^K \exp\{g_k(\mathbf{x})\}}$$



**Goal:** For input features  $\mathbf{x}$ ,  
predict membership in one  
of  $K$  classes:  $C_1, \dots, C_K$

**Training:** Minimize with respect to  $\theta \dots$

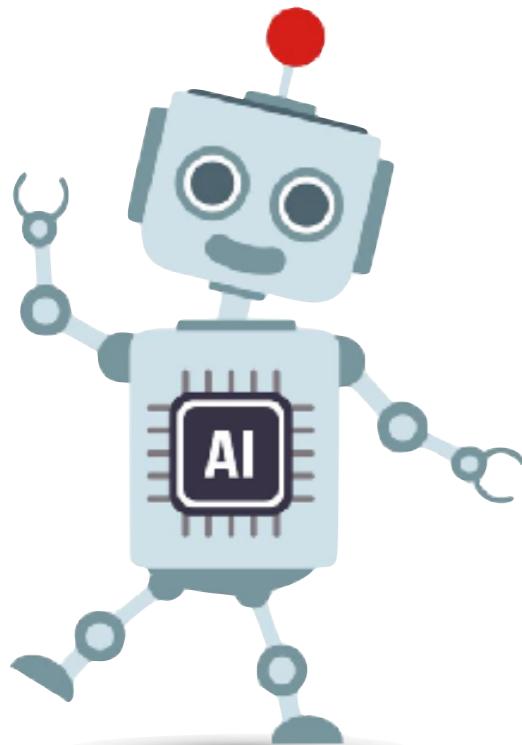


$$\ell(\theta; \mathbf{X}, \mathbf{y}) = -\log \left\{ \prod_n P(C_{y_n} | \mathbf{x}_n) \right\}$$
$$= -\sum_n \log \frac{\exp\{g_{y_n}(\mathbf{x}_n)\}}{\sum_{k=1}^K \exp\{g_k(\mathbf{x}_n)\}}$$

Classifier

**Goal:** For input features  $\mathbf{x}$ ,  
predict membership in one  
of  $K$  classes:  $C_1, \dots, C_K$

**Training:** Minimize with respect to  $\theta \dots$

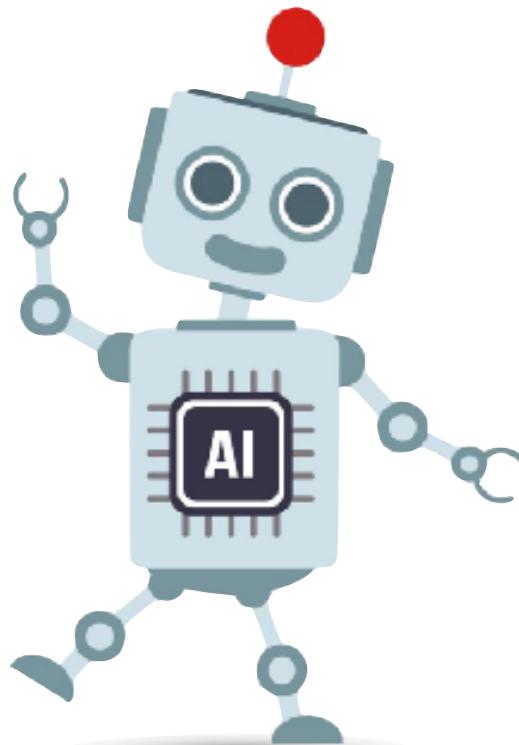


$$\ell(\theta; \mathbf{X}, \mathbf{y}) = -\log \left\{ \prod_n P(C_{y_n} | \mathbf{x}_n) \right\}$$
$$= -\sum_n \log \frac{\exp\{g_{y_n}(\mathbf{x}_n)\}}{\sum_{k=1}^K \exp\{g_k(\mathbf{x}_n)\}}$$

Classifier

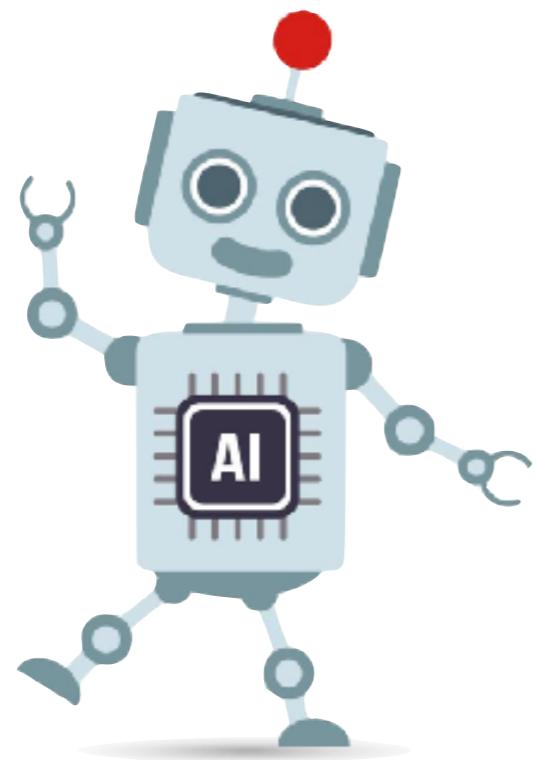
Goal: For input features  $\mathbf{x}$ ,  
predict membership in one  
of  $K$  classes:  $C_1, \dots, C_K$

Training: Minimize with respect to  $\theta \dots$

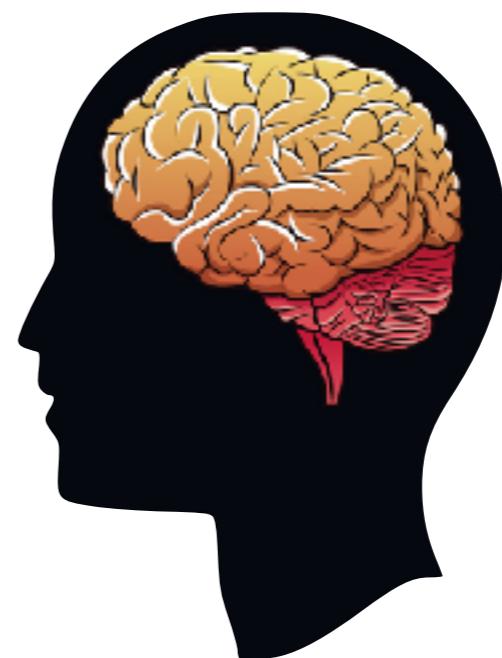
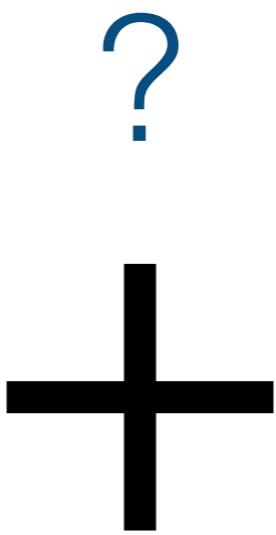


$$\ell(\theta; \mathbf{X}, \mathbf{y}) = -\log \left\{ \prod_n P(C_{y_n} | \mathbf{x}_n) \right\}$$
$$= - \sum_n \log \frac{\exp\{g_{y_n}(\mathbf{x}_n)\}}{\sum_{k=1}^K \exp\{g_k(\mathbf{x}_n)\}}$$

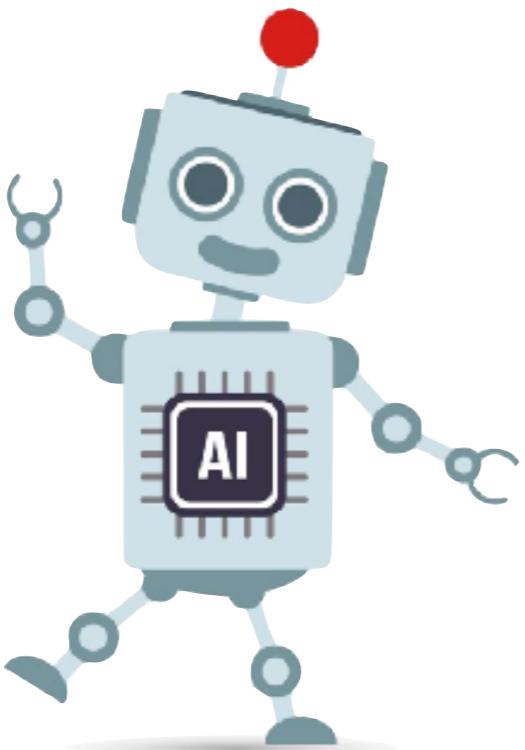
Classifier



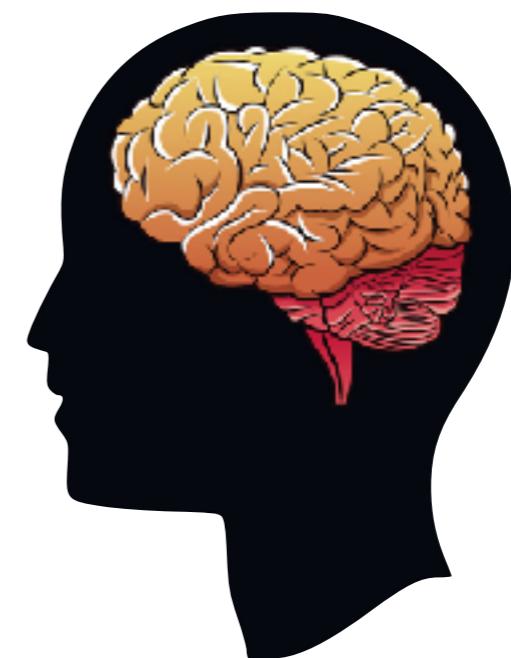
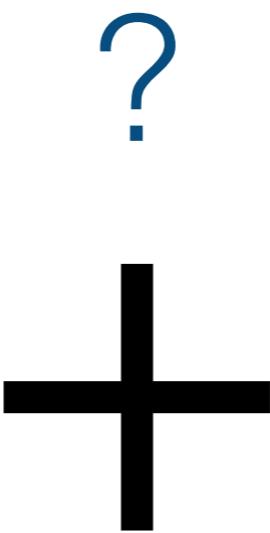
Classifier



Expert



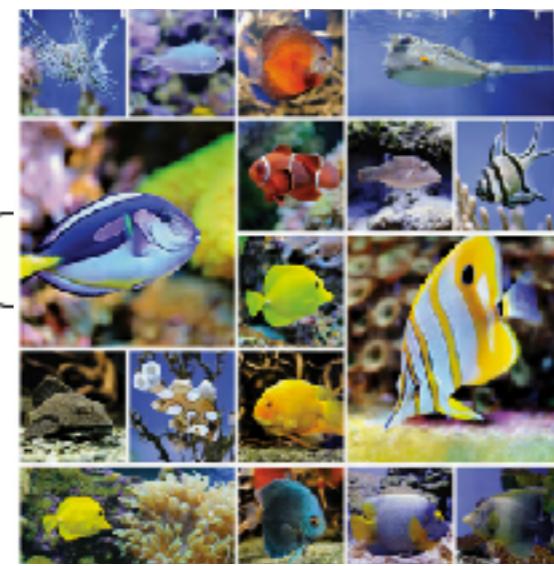
Classifier

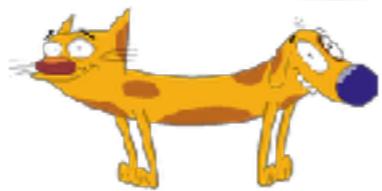
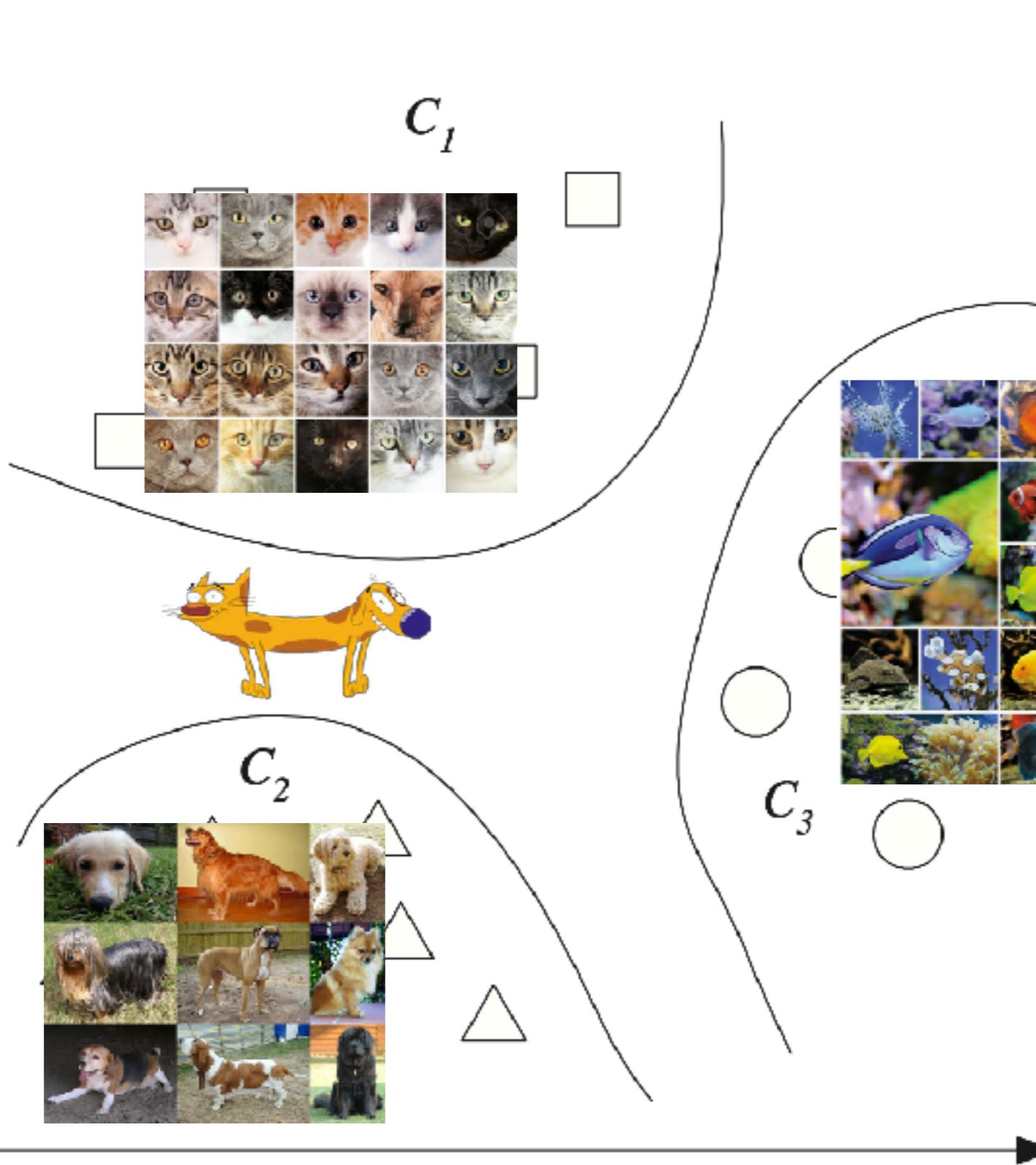


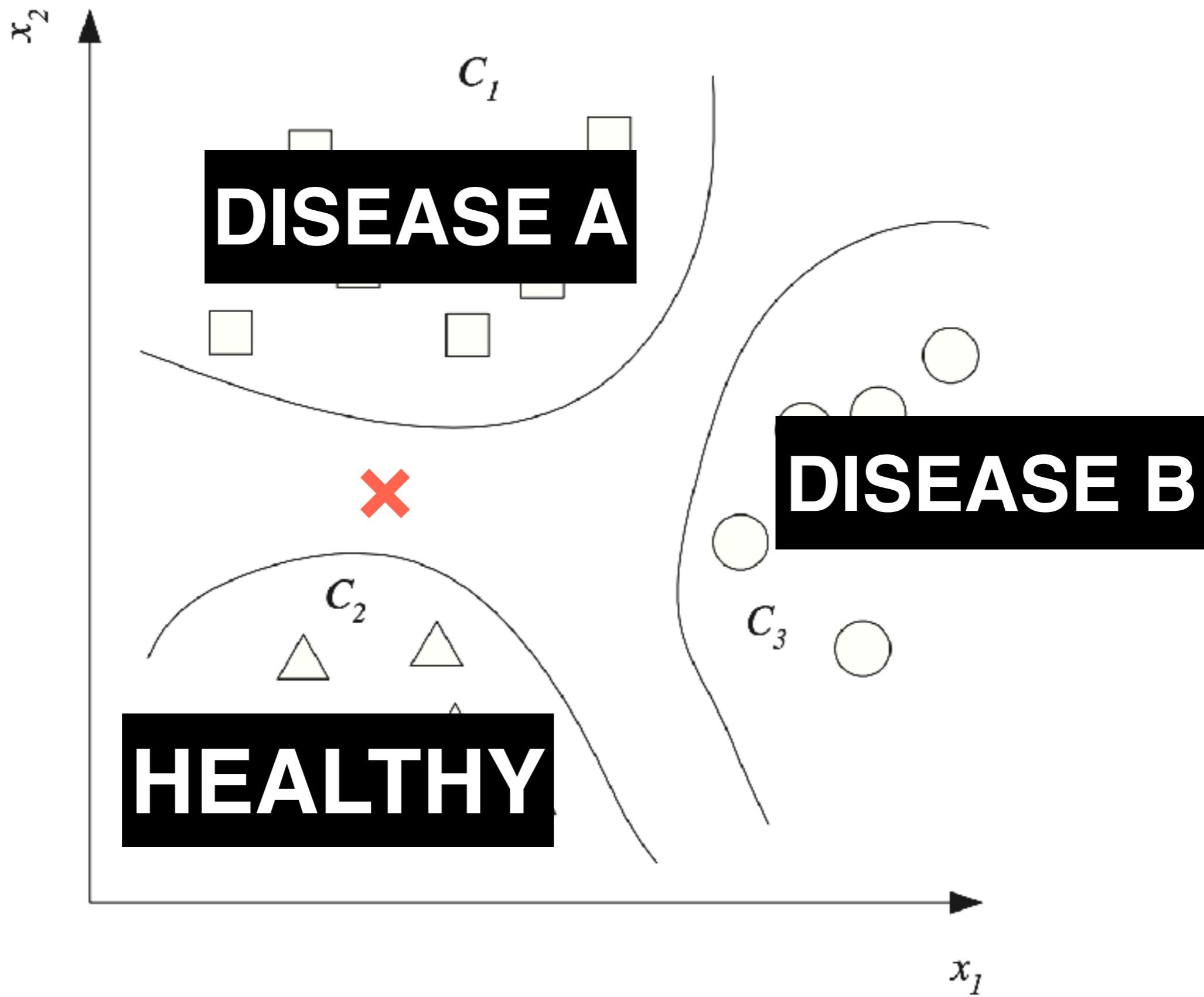
Expert

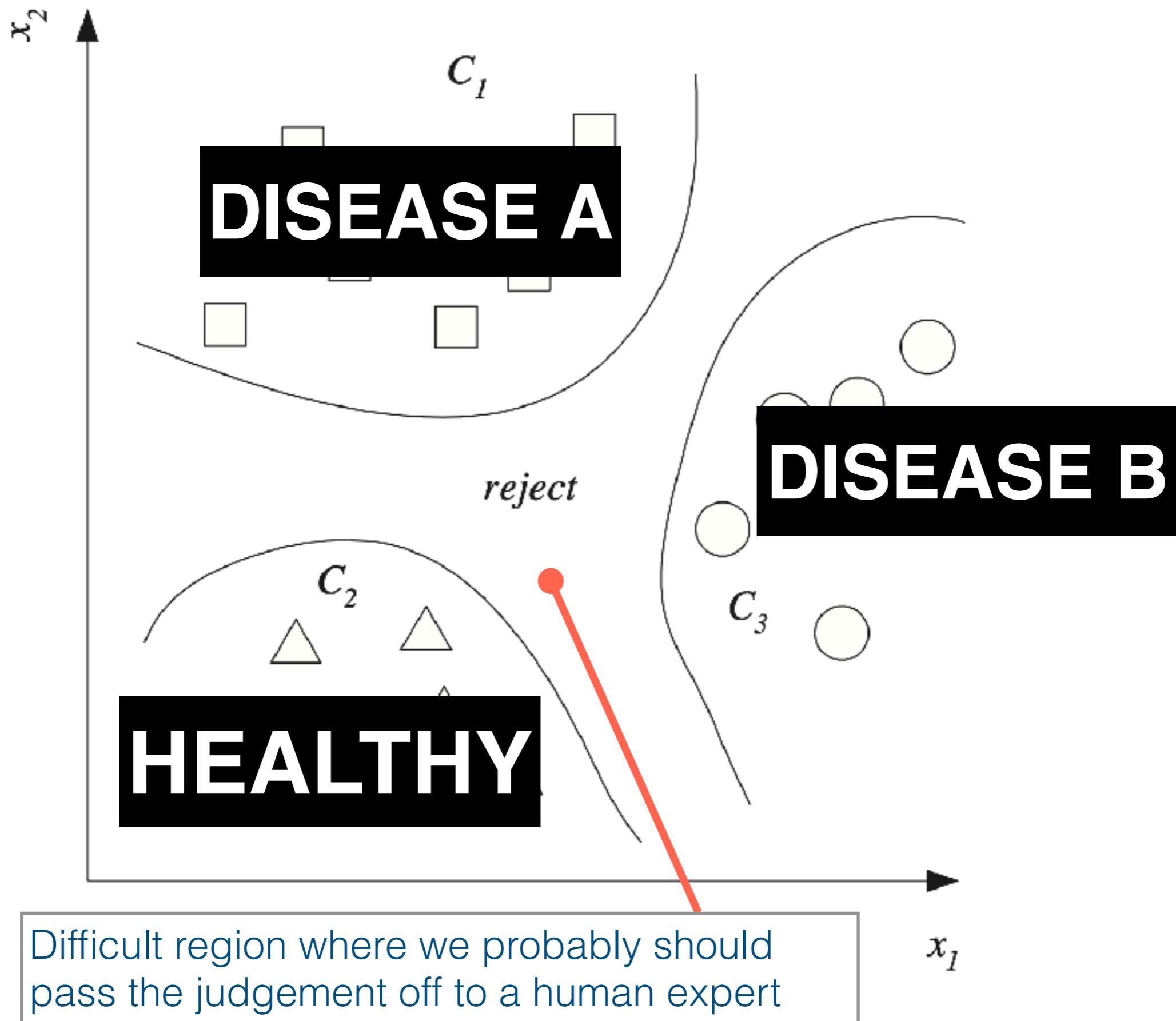
# Warm Up:

## Classification with a Rejection Option

$x_2$  $C_1$  $C_2$  $C_3$  $x_1$  $\rightarrow$

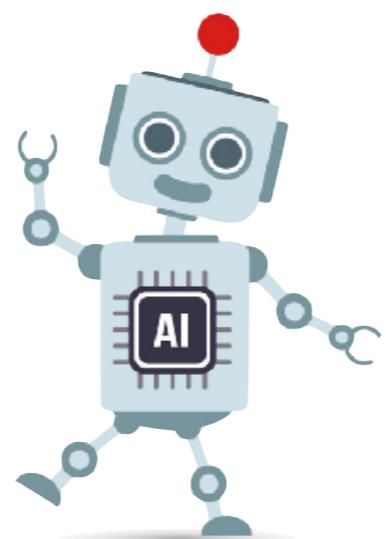
$x_2$  $C_1$  $C_2$  $C_3$  $x_1$ 







X

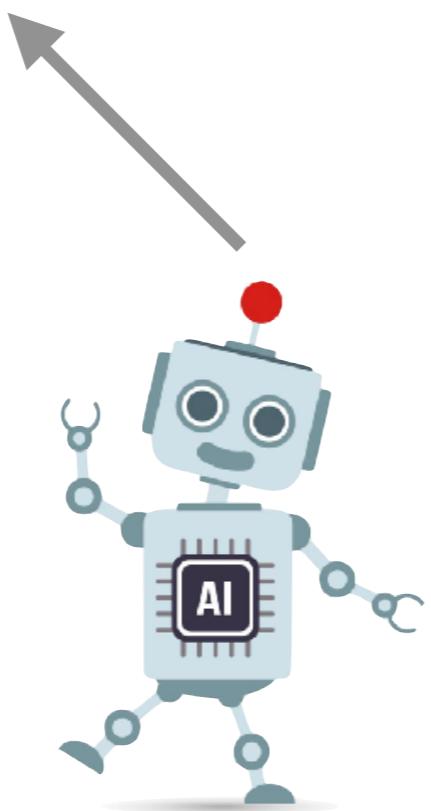


Classifier

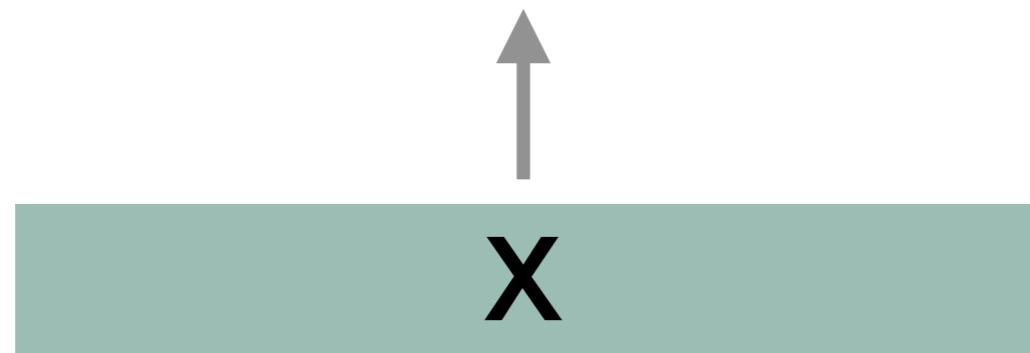


X

make  
prediction

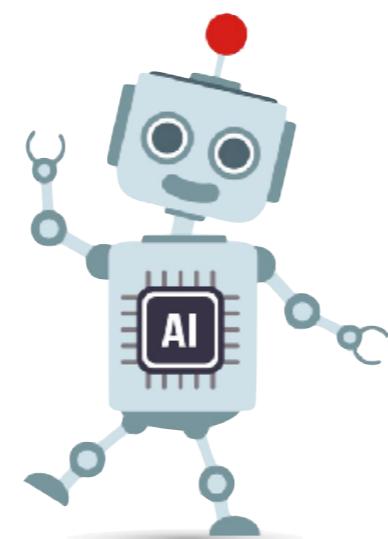


Classifier



make  
prediction

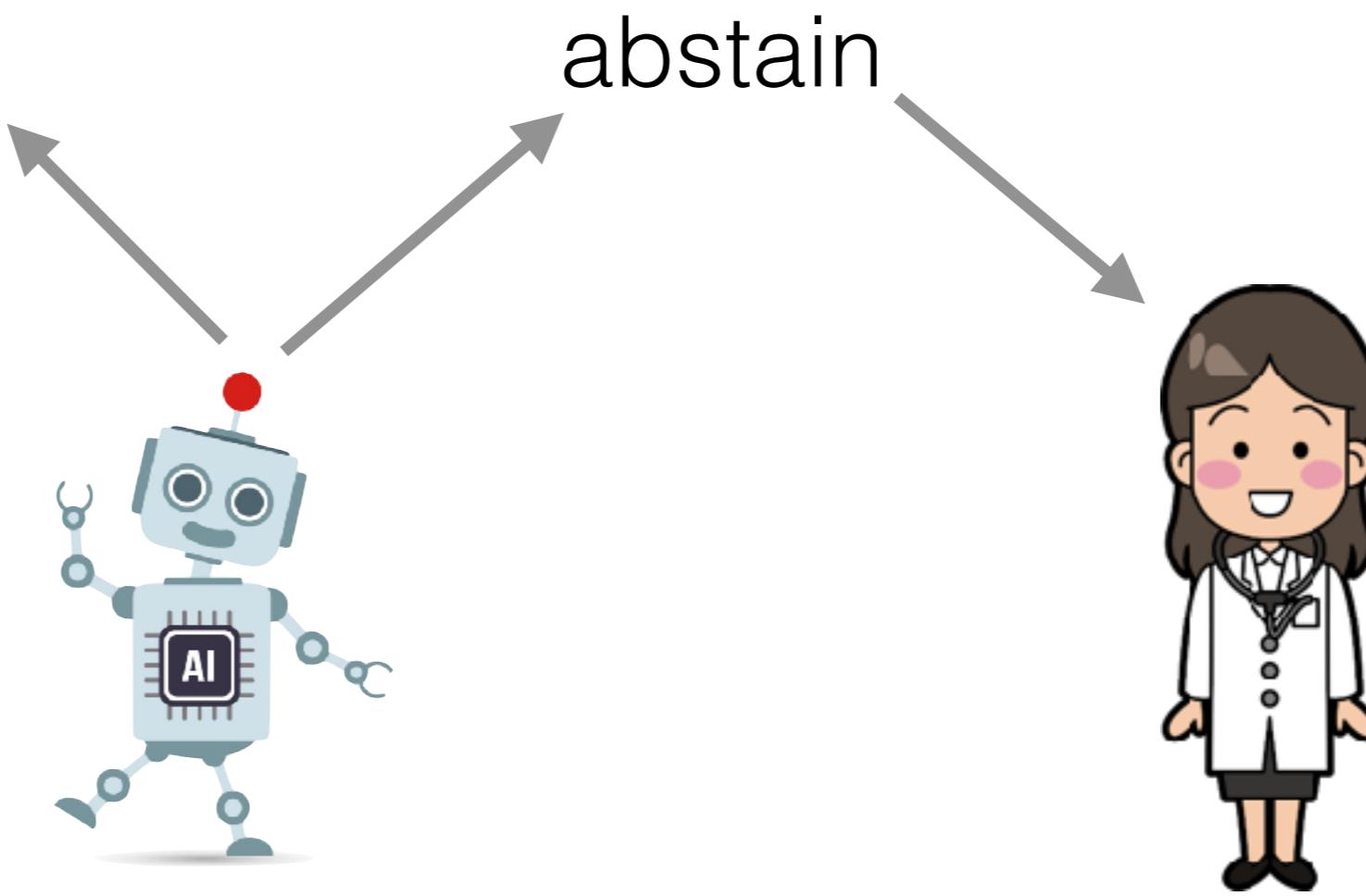
abstain



Classifier

X

make  
prediction

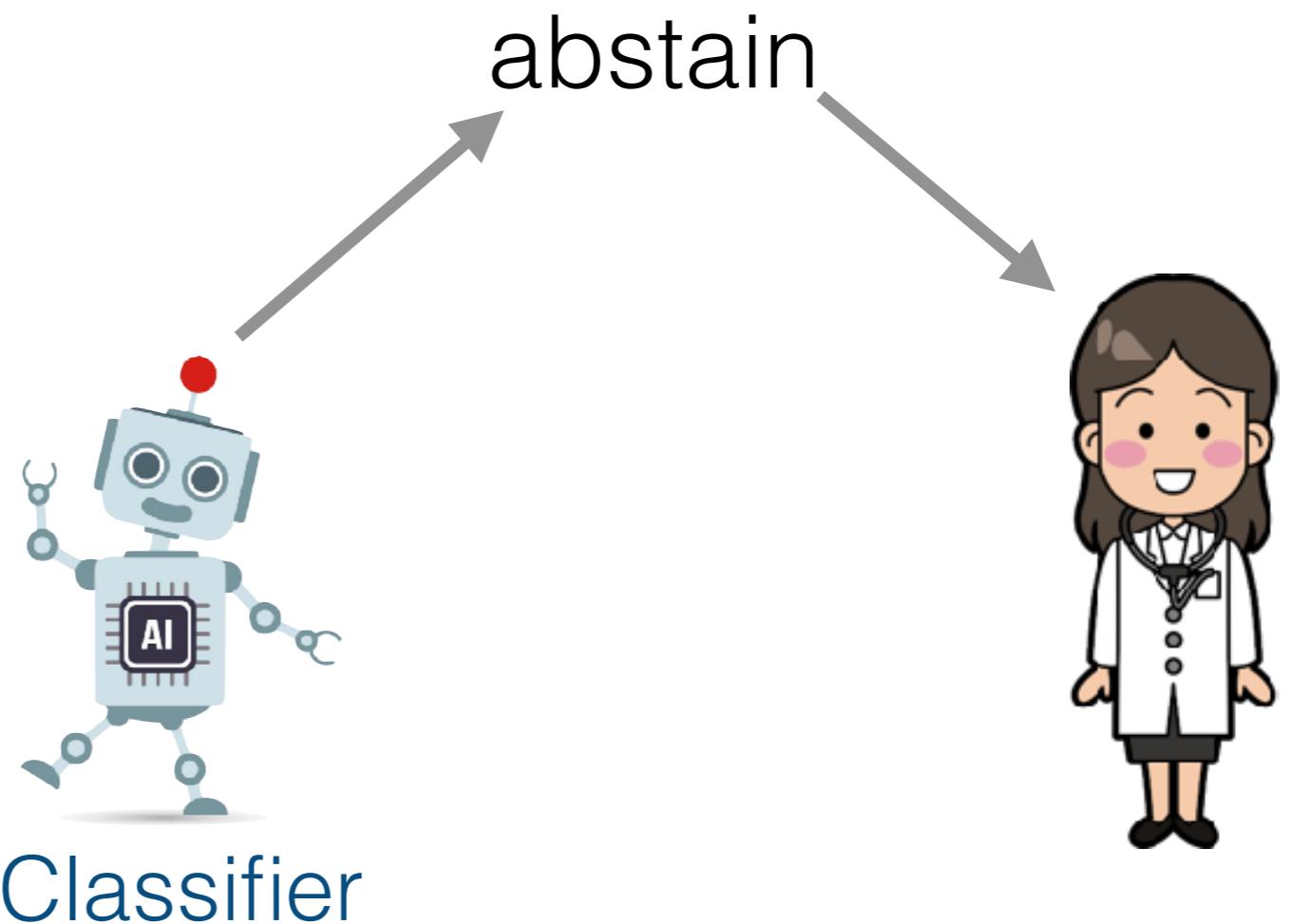


Classifier

X

Confidence-Based Rejection: Abstain if the model is unconfident in its prediction:

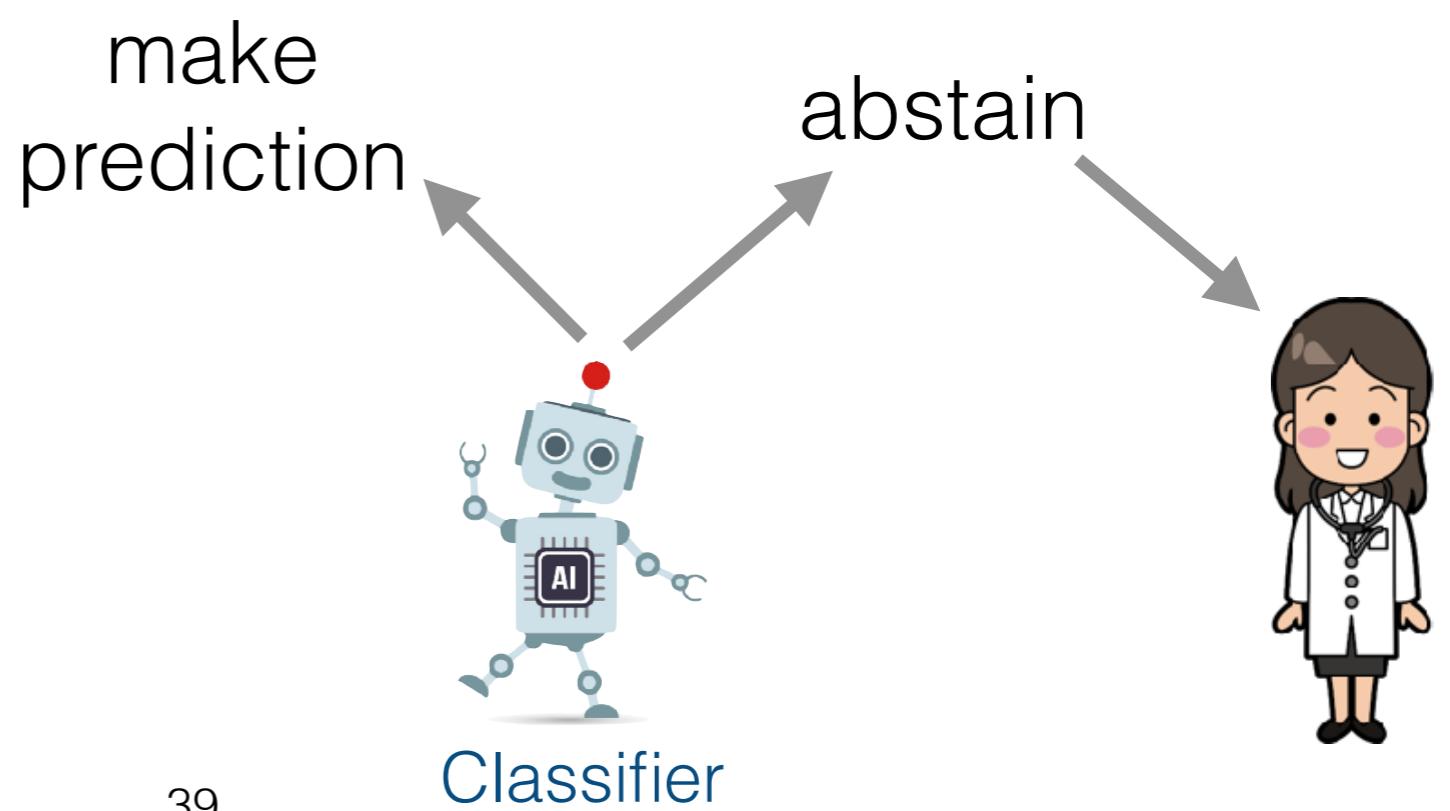
$$\max_k P(C_k | x_n) < \tau$$



# Learning to Defer to an Expert

# Issue with Rejection Option

- ⊗ The human is the “backup plan”: catching the inputs that the classifier is unsure about.



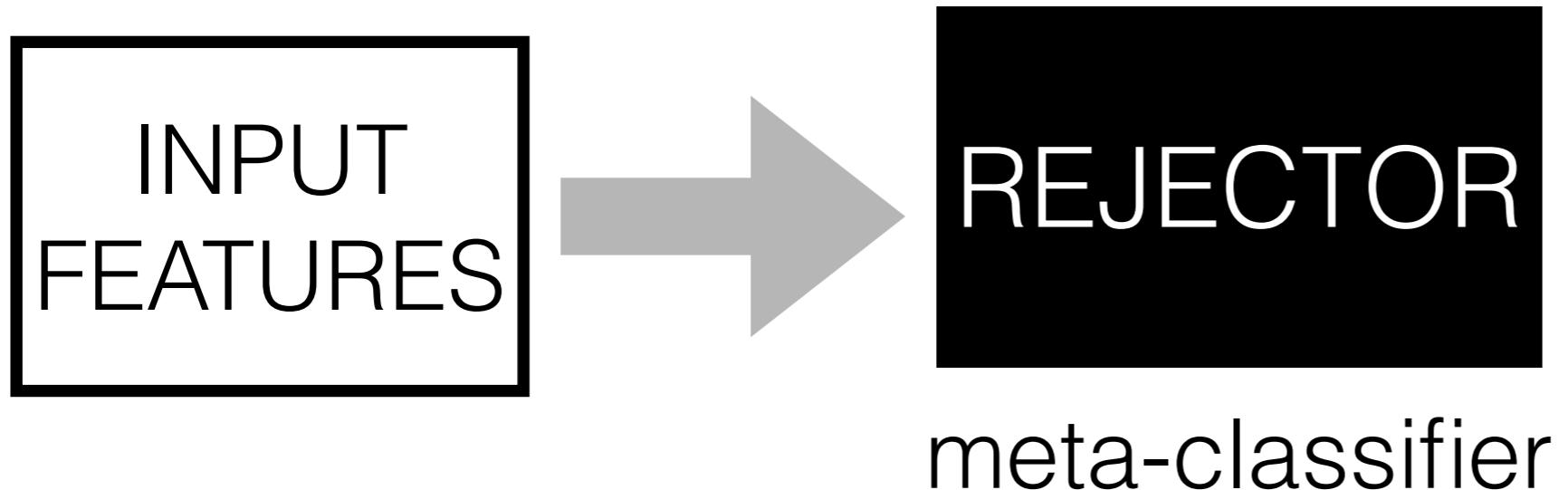
# Issue with Rejection Option

- ⊗ The human is the “backup plan”: catching the inputs that the classifier is unsure about.

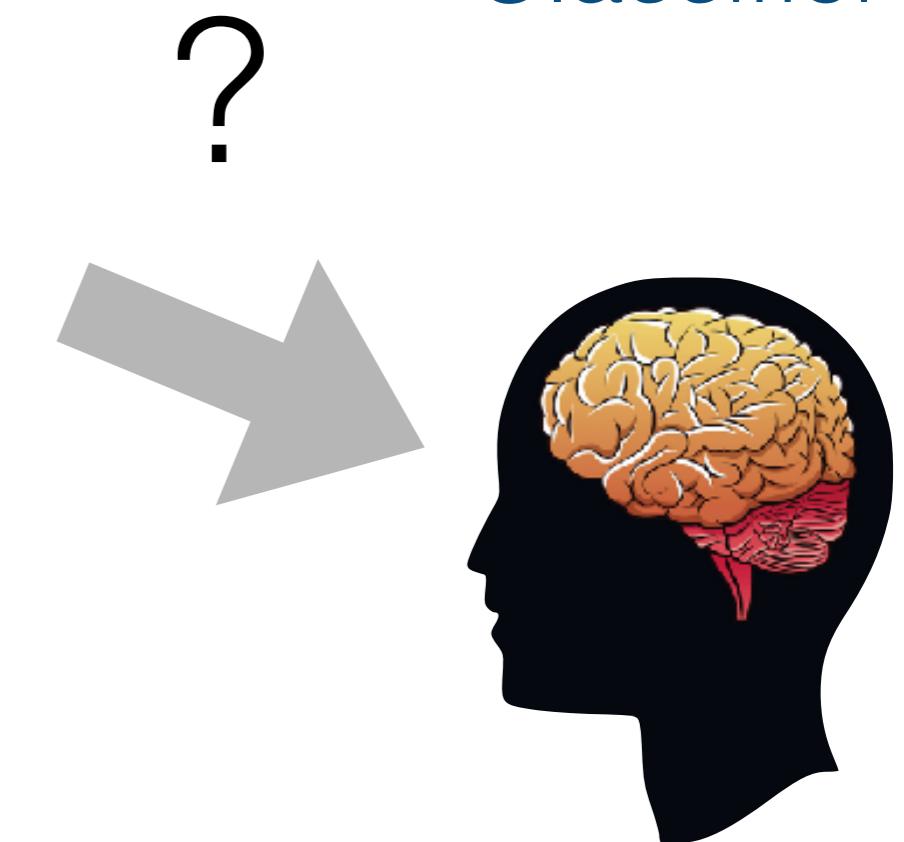
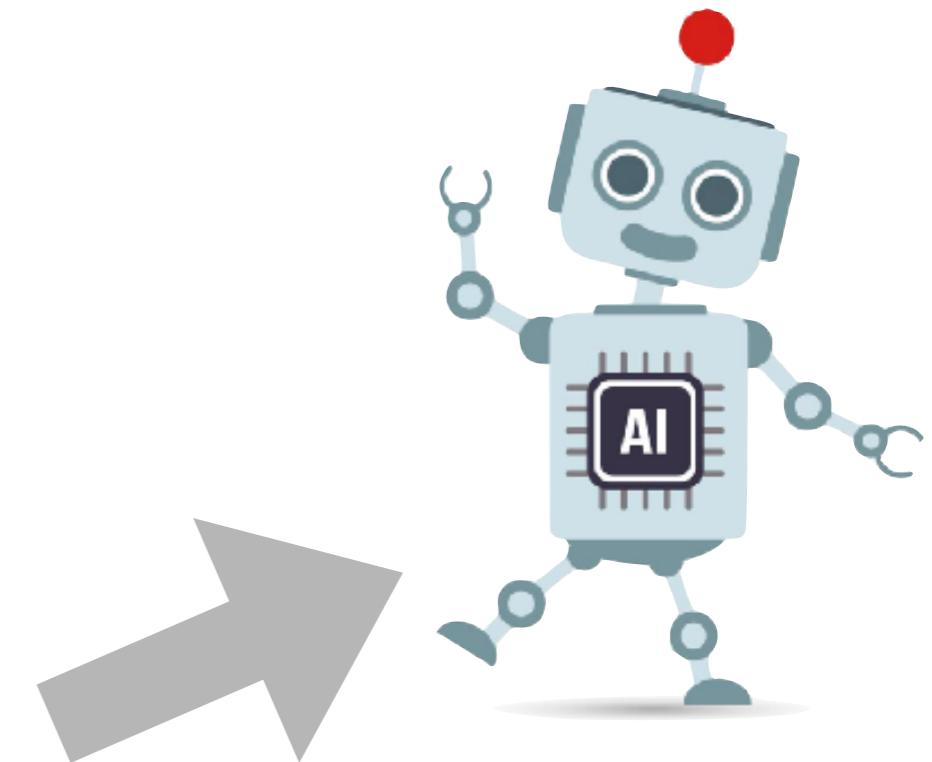
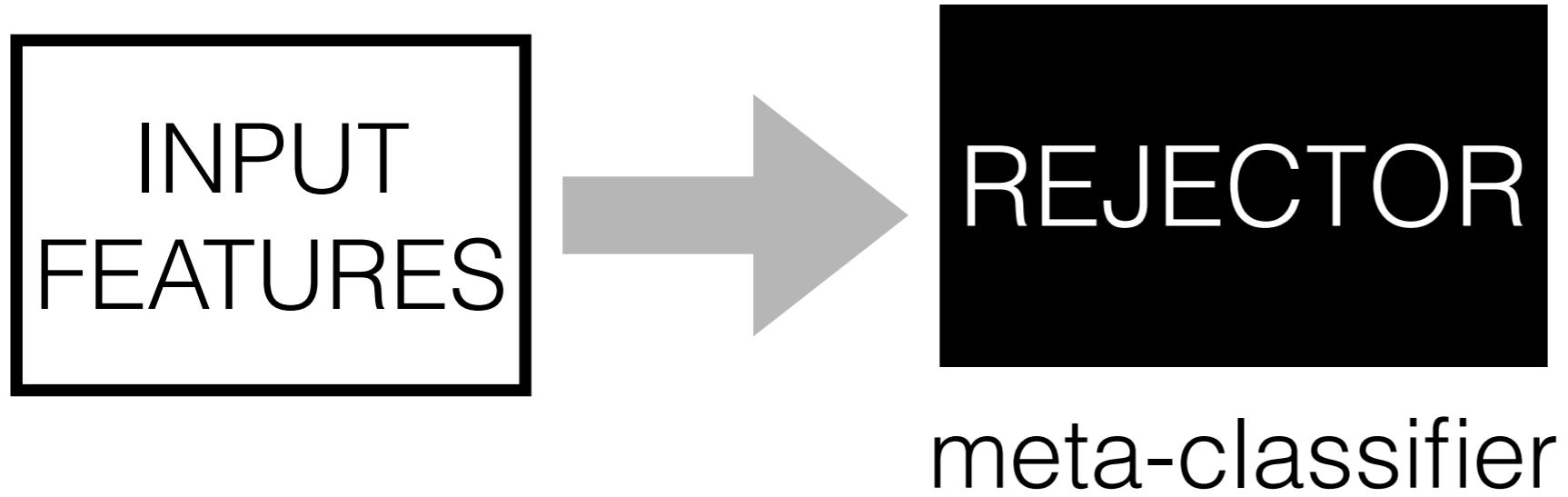
## Better Formulation?

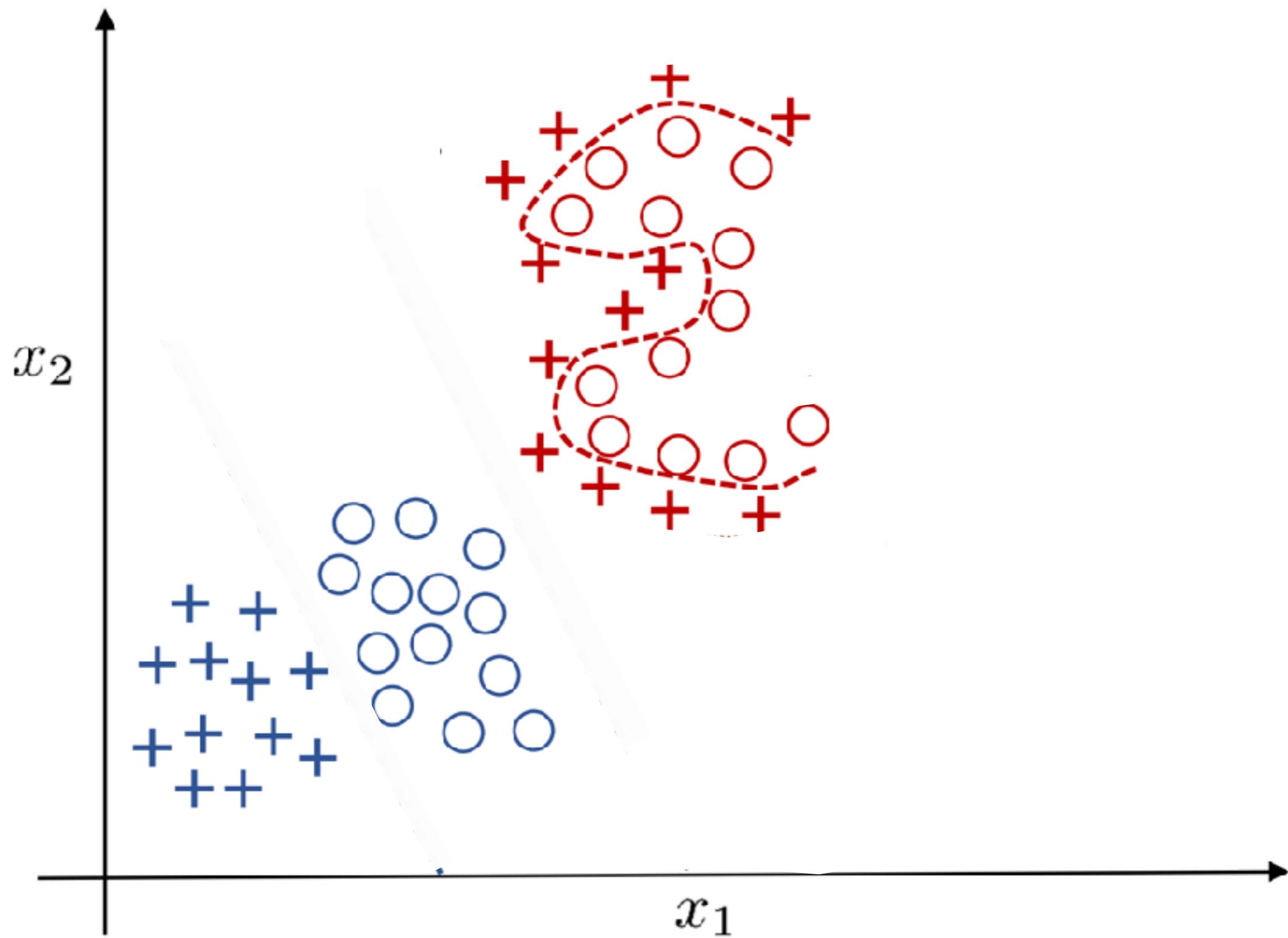
- ⊗ Model what the human knows as well, so we can enable *collaboration*.

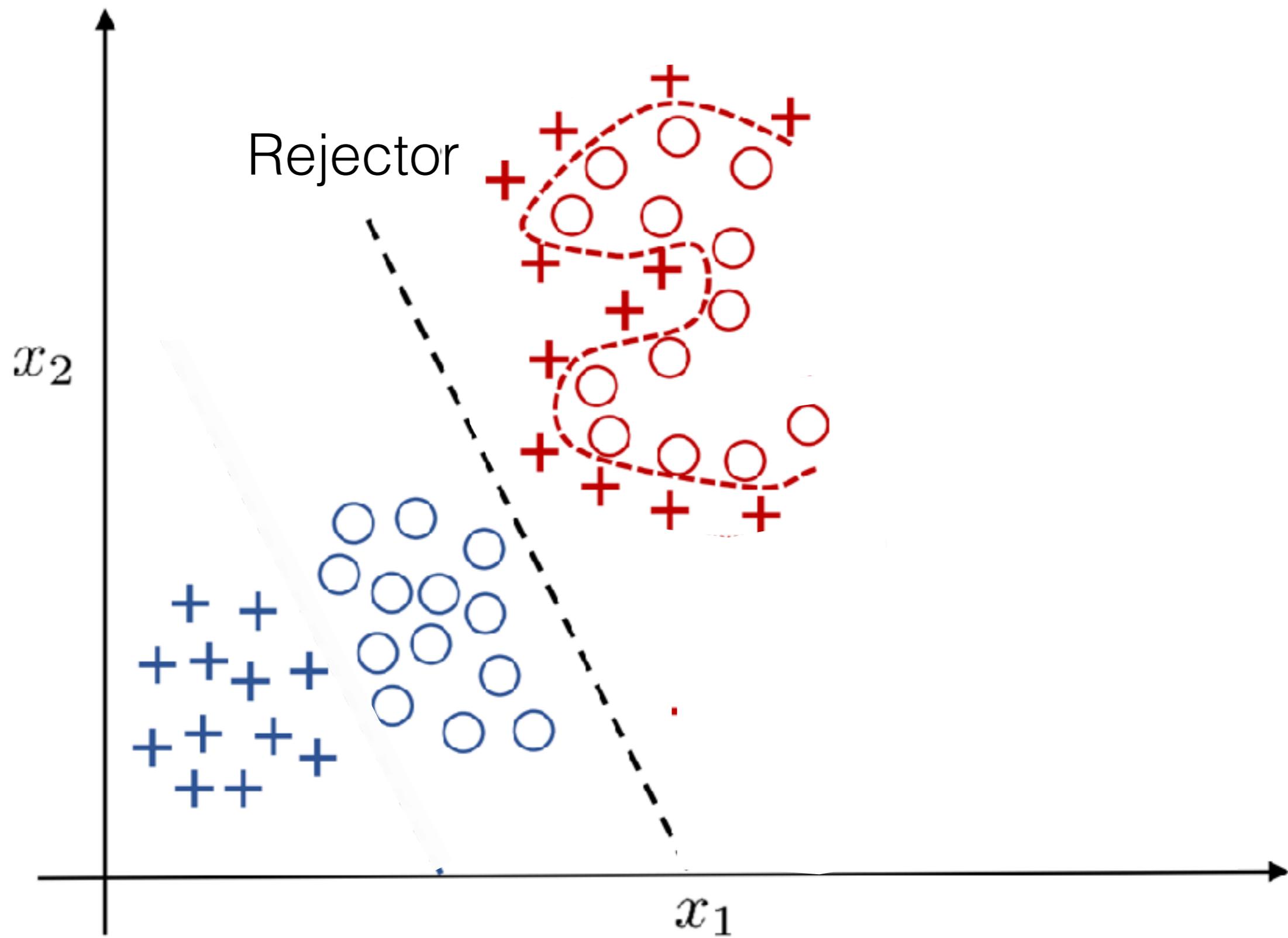
# Learning to Defer

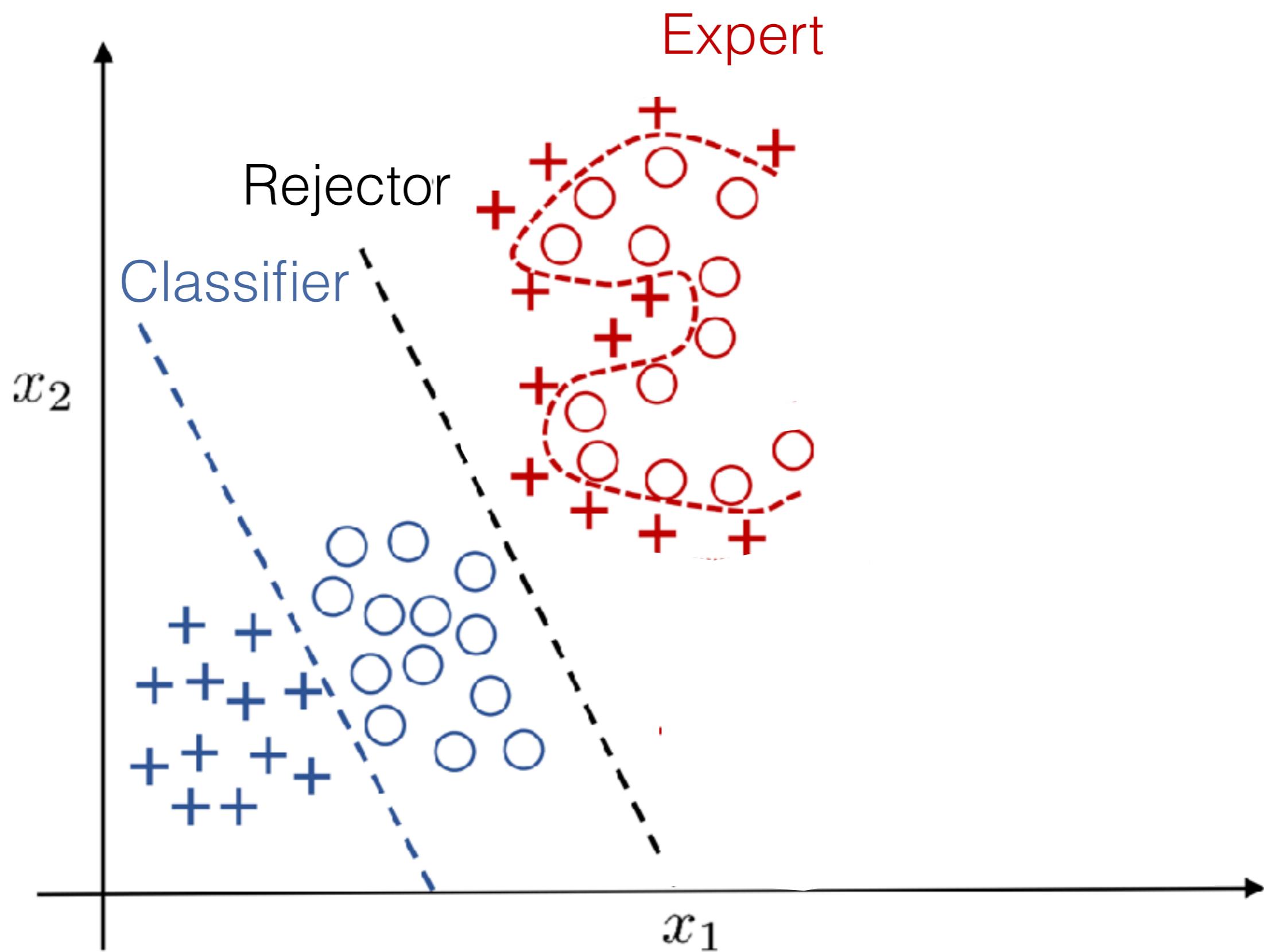


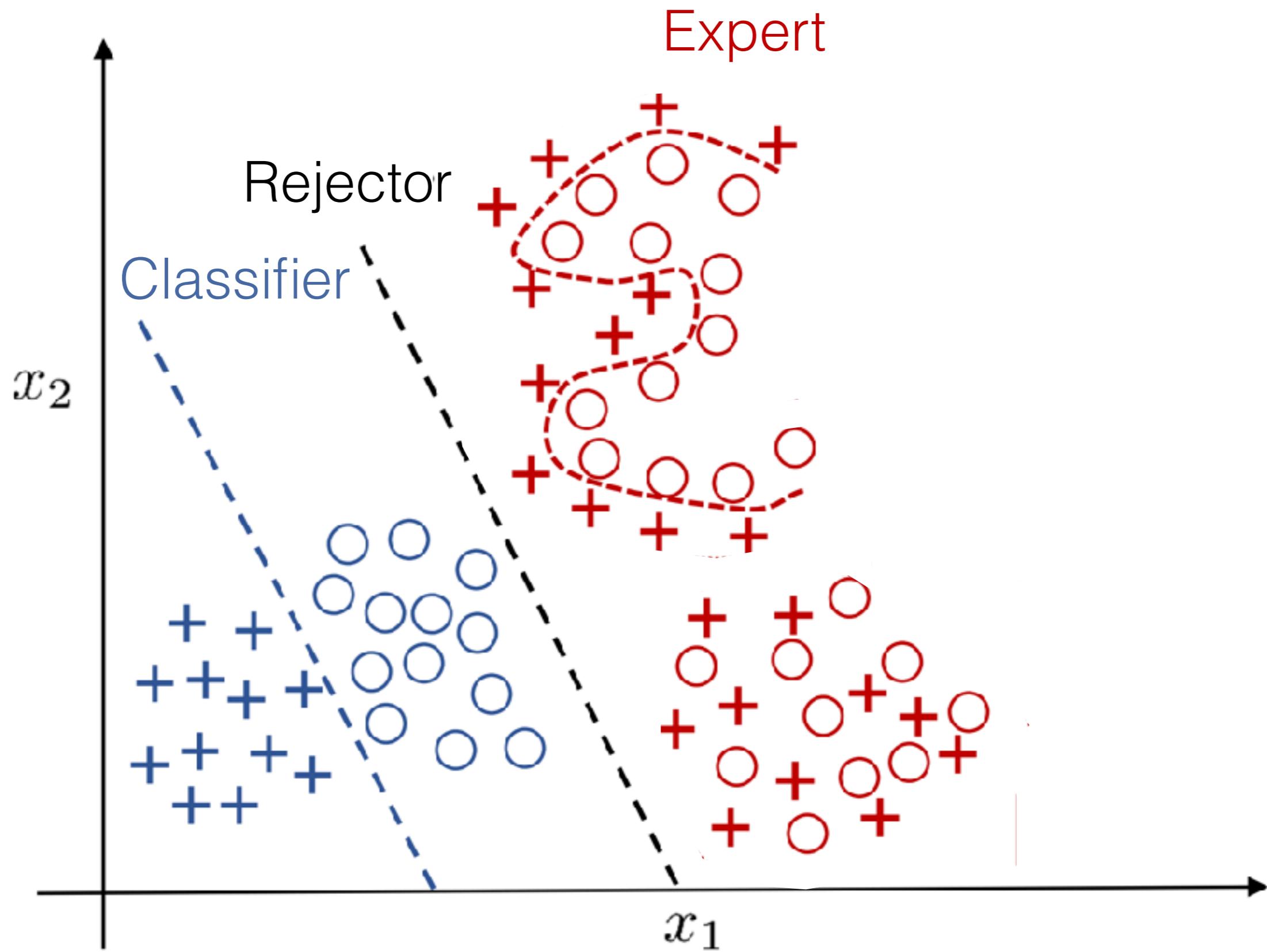
# Learning to Defer



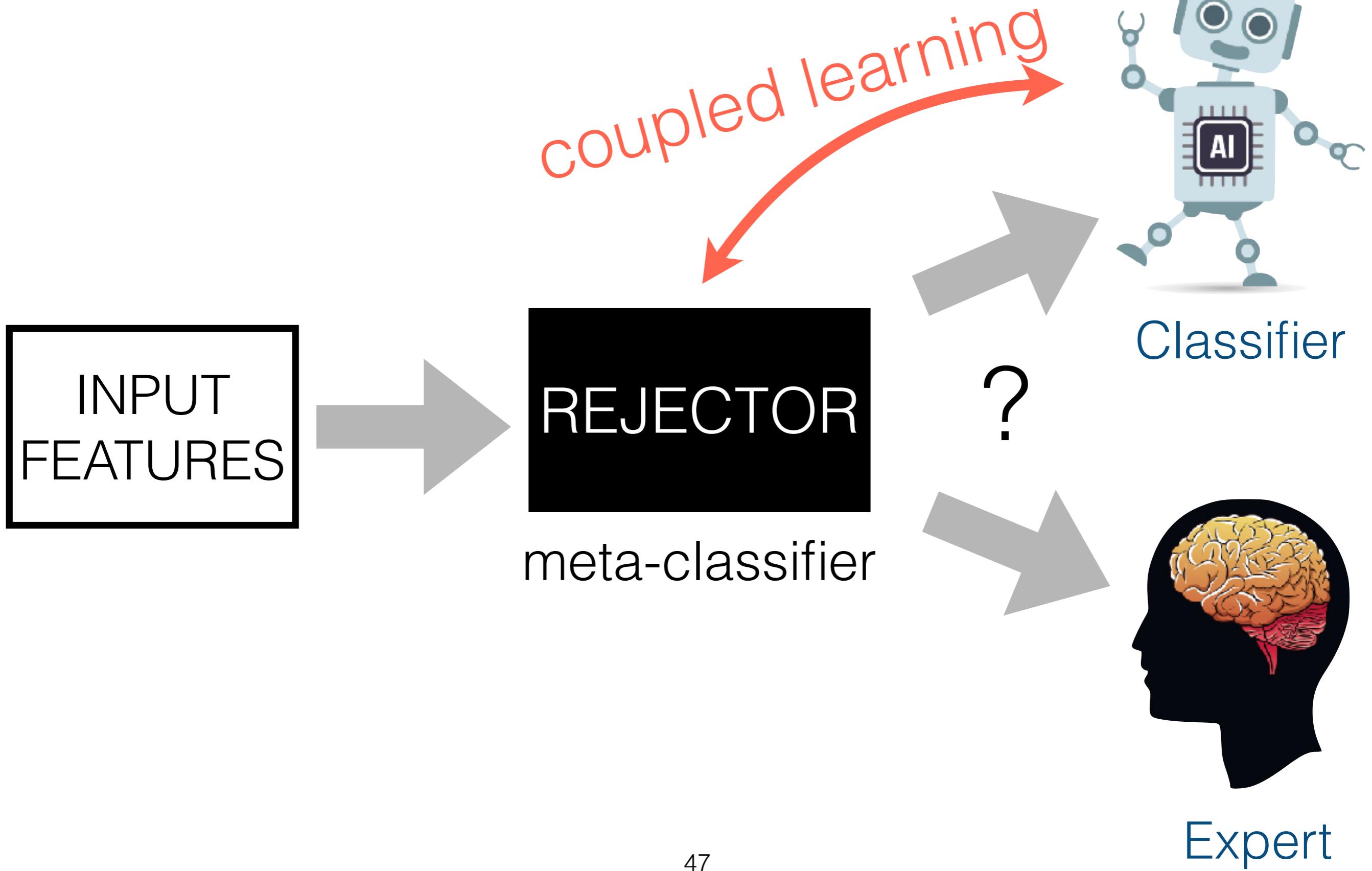








# Learning to Defer



# Learning to Defer

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

*expert predictions*



# Learning to Defer

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

expert predictions

Models:

$r(\mathbf{x})$	$h(\mathbf{x})$
Rejector	Classifier

# Learning to Defer

0-1 Loss:

$$\ell(r, h; \mathcal{D}) =$$

$$\sum_n (1 - r(x_n)) \mathbb{I}[h(x_n) \neq y_n] + r(x_n) \mathbb{I}[m_n \neq y_n]$$


classifier loss      expert loss

# Learning to Defer

0-1 Loss:

$$\ell(r, h; \mathcal{D}) =$$

$$\sum_n (1 - r(x_n)) \mathbb{I}[h(x_n) \neq y_n] + r(x_n) \mathbb{I}[m_n \neq y_n]$$


classifier loss      expert loss

# Learning to Defer

0-1 Loss:

$$\ell(r, h; \mathcal{D}) =$$

$$\sum_n (1 - r(x_n)) \mathbb{I}[h(x_n) \neq y_n] + r(x_n) \mathbb{I}[m_n \neq y_n]$$

classifier loss

expert loss

# Learning to Defer

Formerly open problem:

what's a consistent surrogate  
for the 0-1 loss?

# Learning to Defer

Formerly open problem:

what's a consistent surrogate  
for the 0-1 loss?

$$h^*(x) = \operatorname{argmax}_y \mathbb{P}(y = y | x)$$

# Learning to Defer

Formerly open problem:

what's a consistent surrogate  
for the 0-1 loss?

$$h^*(x) = \operatorname{argmax}_y \mathbb{P}(y = y | x)$$

$$r^*(x) = \mathbb{I} [\mathbb{P}(m = y | x) \geq \max_y \mathbb{P}(y = y | x)]$$

probability that the expert is correct

# Consistent Estimators for Learning to Defer to an Expert

Hussein Mozannar \*      David Sontag †

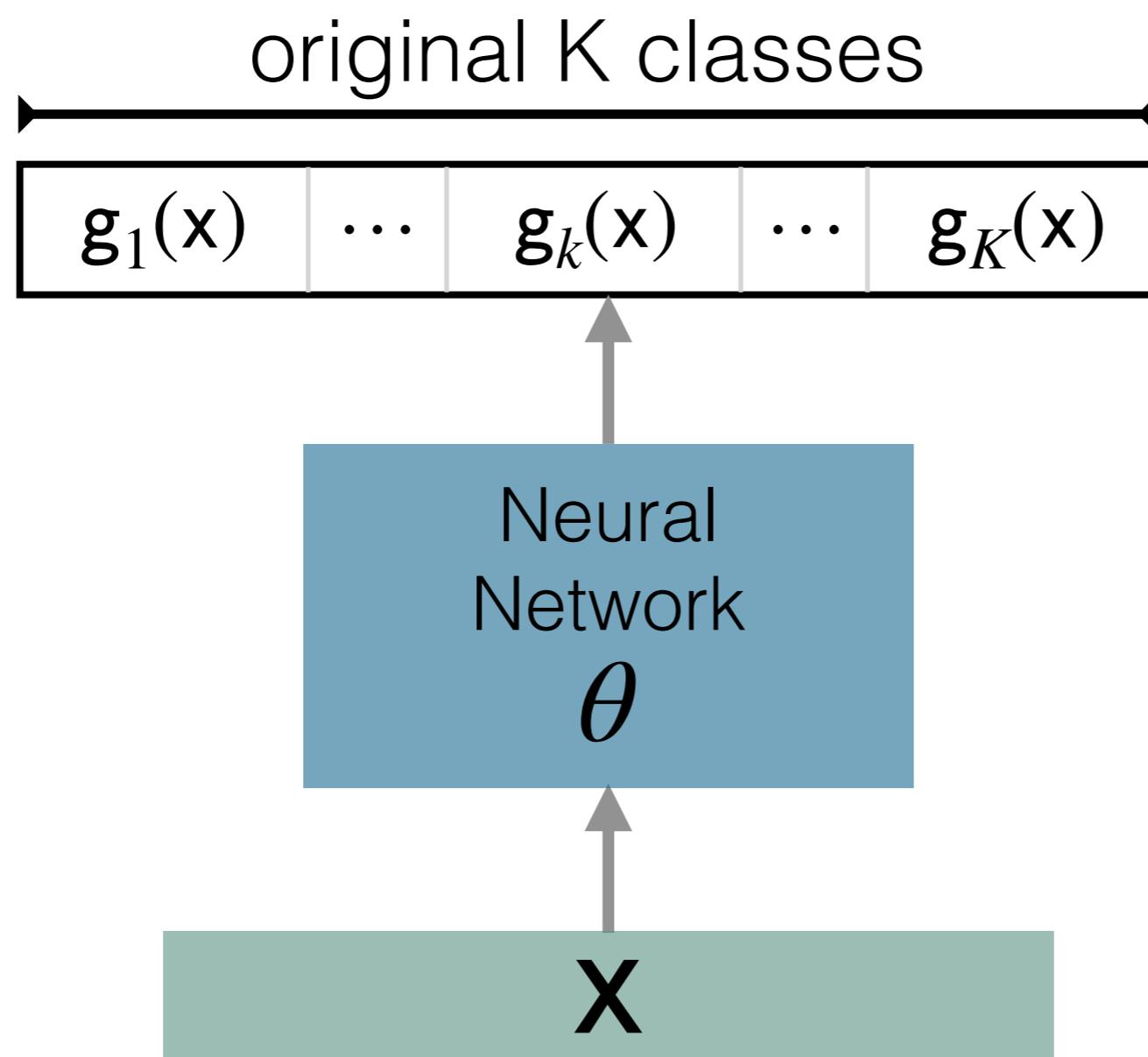
## Abstract

Learning algorithms are often used in conjunction with expert decision makers in practical scenarios, however this fact is largely ignored when designing these algorithms. In this paper we explore how to learn predictors that can either predict or choose to defer the decision to a downstream expert. Given only samples of the expert's decisions, we give a procedure based on learning a classifier and a rejector and analyze it theoretically. Our approach is based on a novel reduction to cost sensitive learning where we give a consistent surrogate loss for cost sensitive learning that generalizes the cross entropy loss. We show the effectiveness of our approach on a variety of experimental tasks.

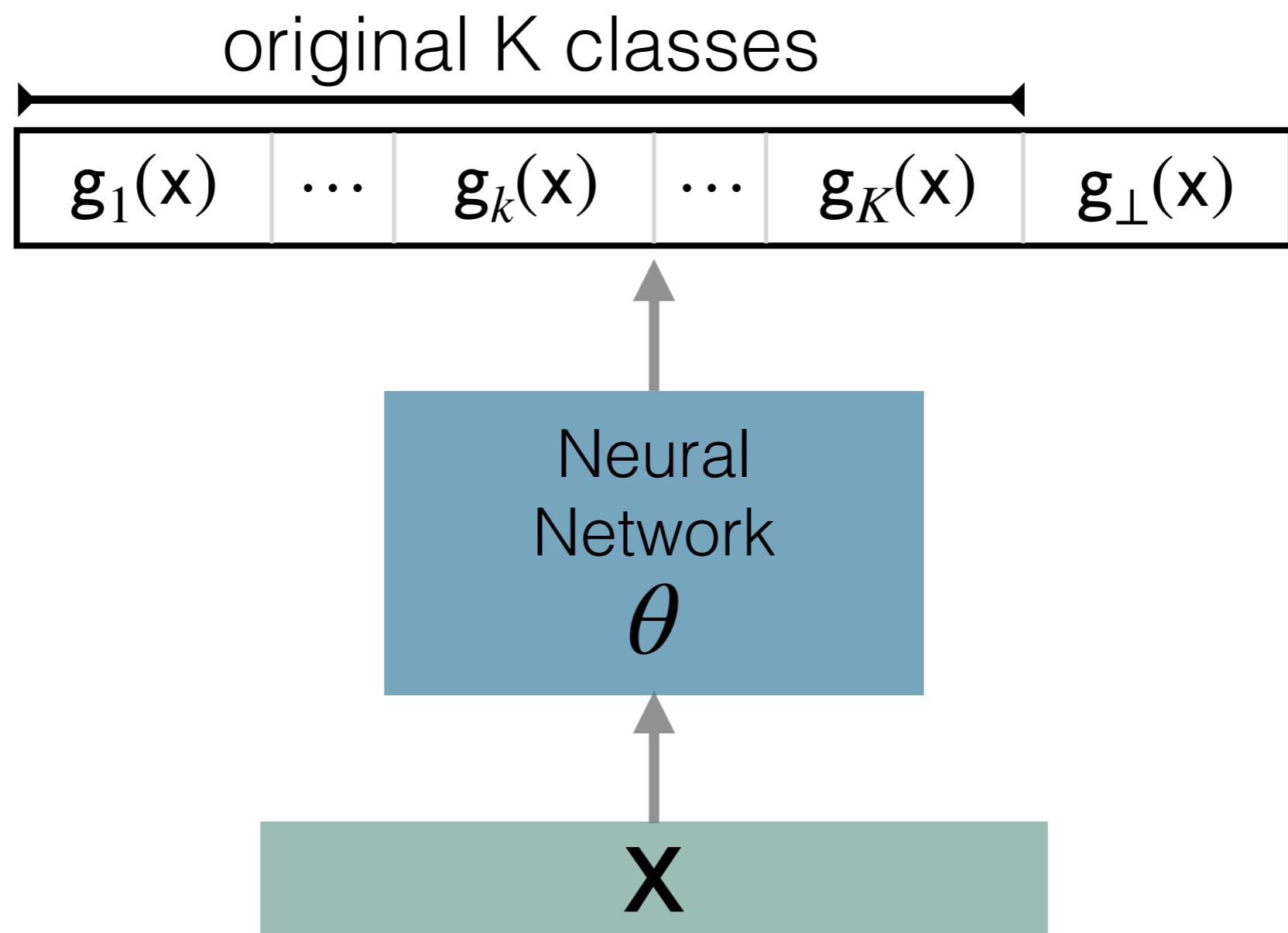
ICML 2020

Solution: Merge the rejector and classifier.

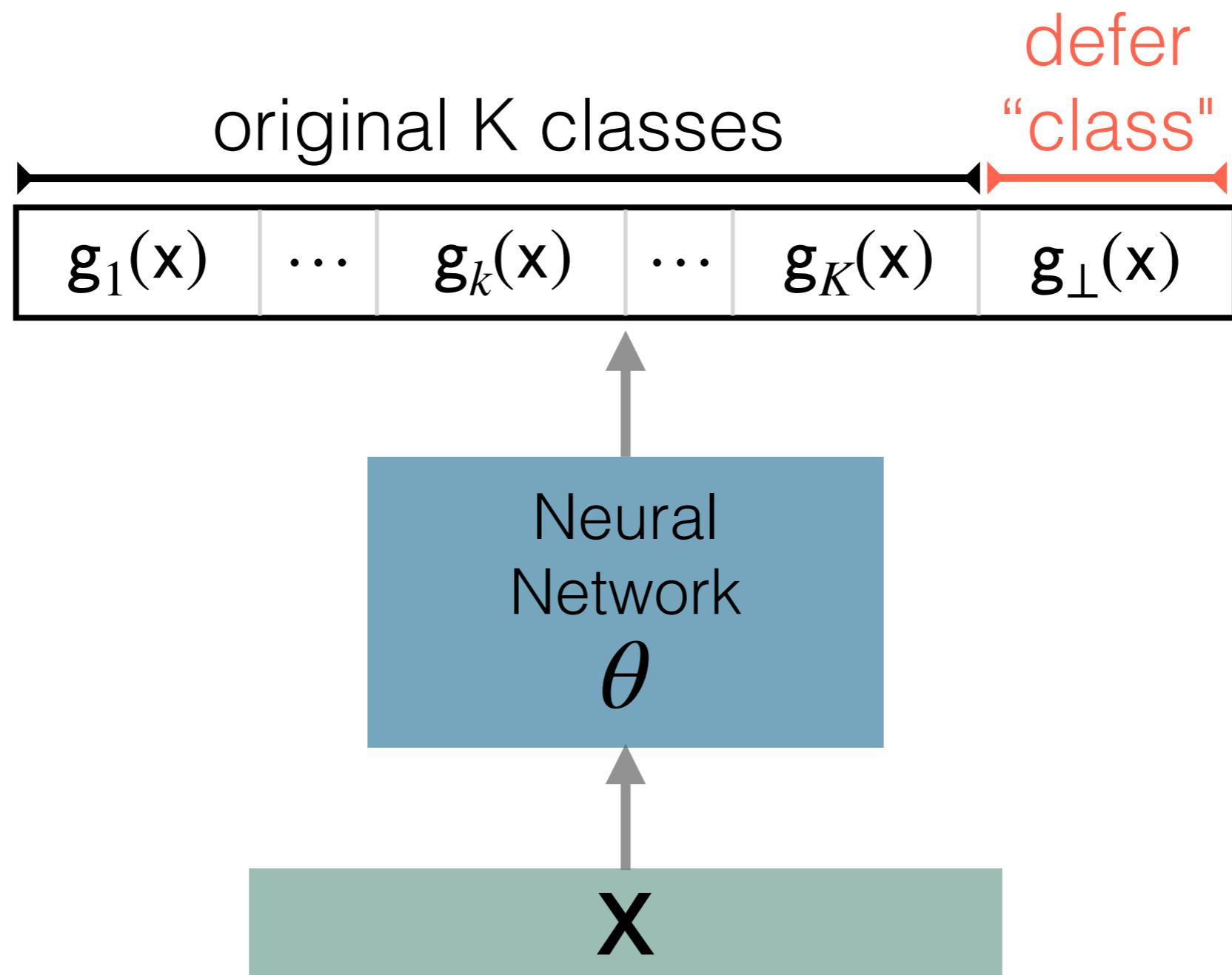
# Solution: Merge the rejector and classifier.



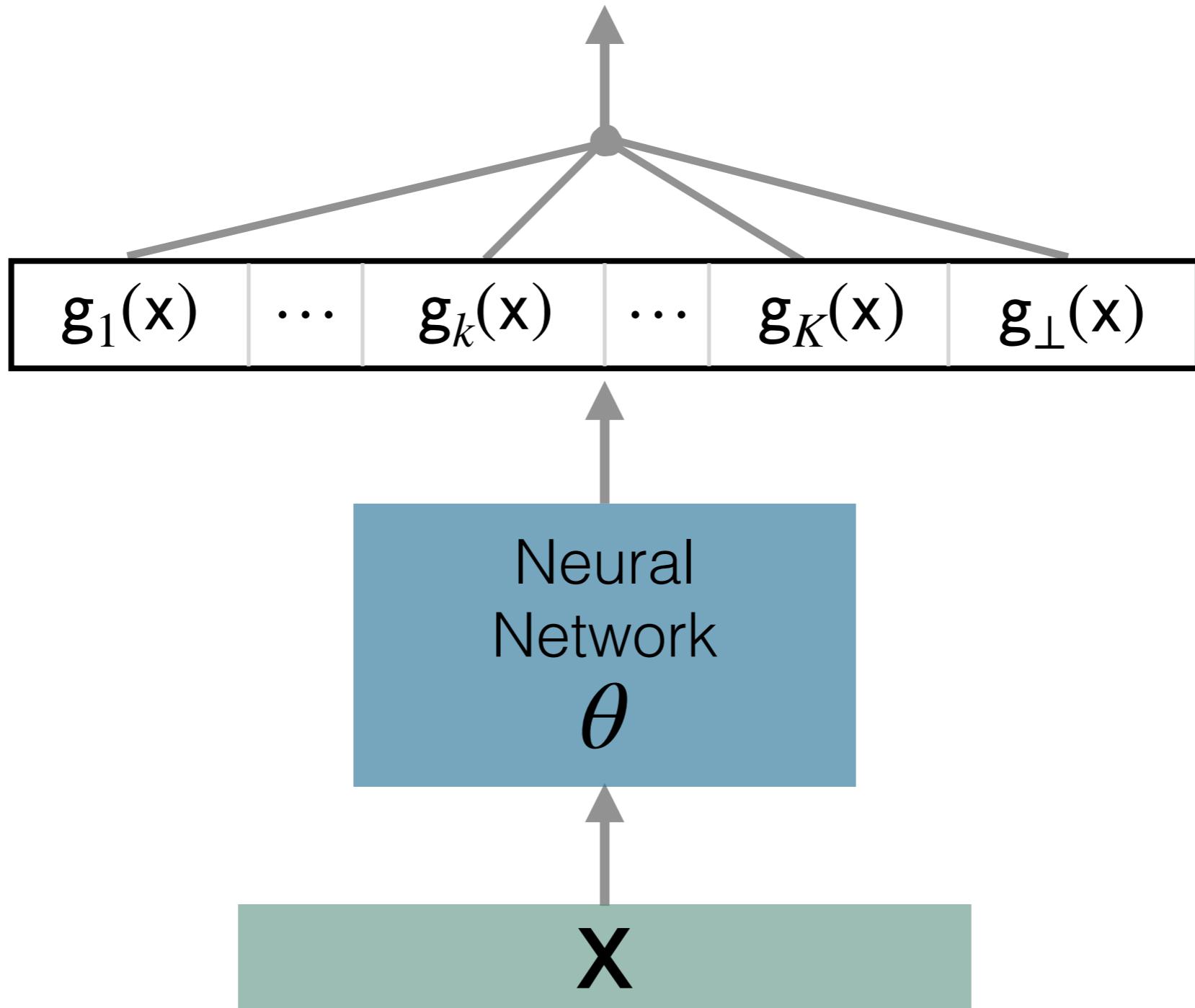
# Solution: Merge the rejector and classifier.



# Solution: Merge the rejector and classifier.

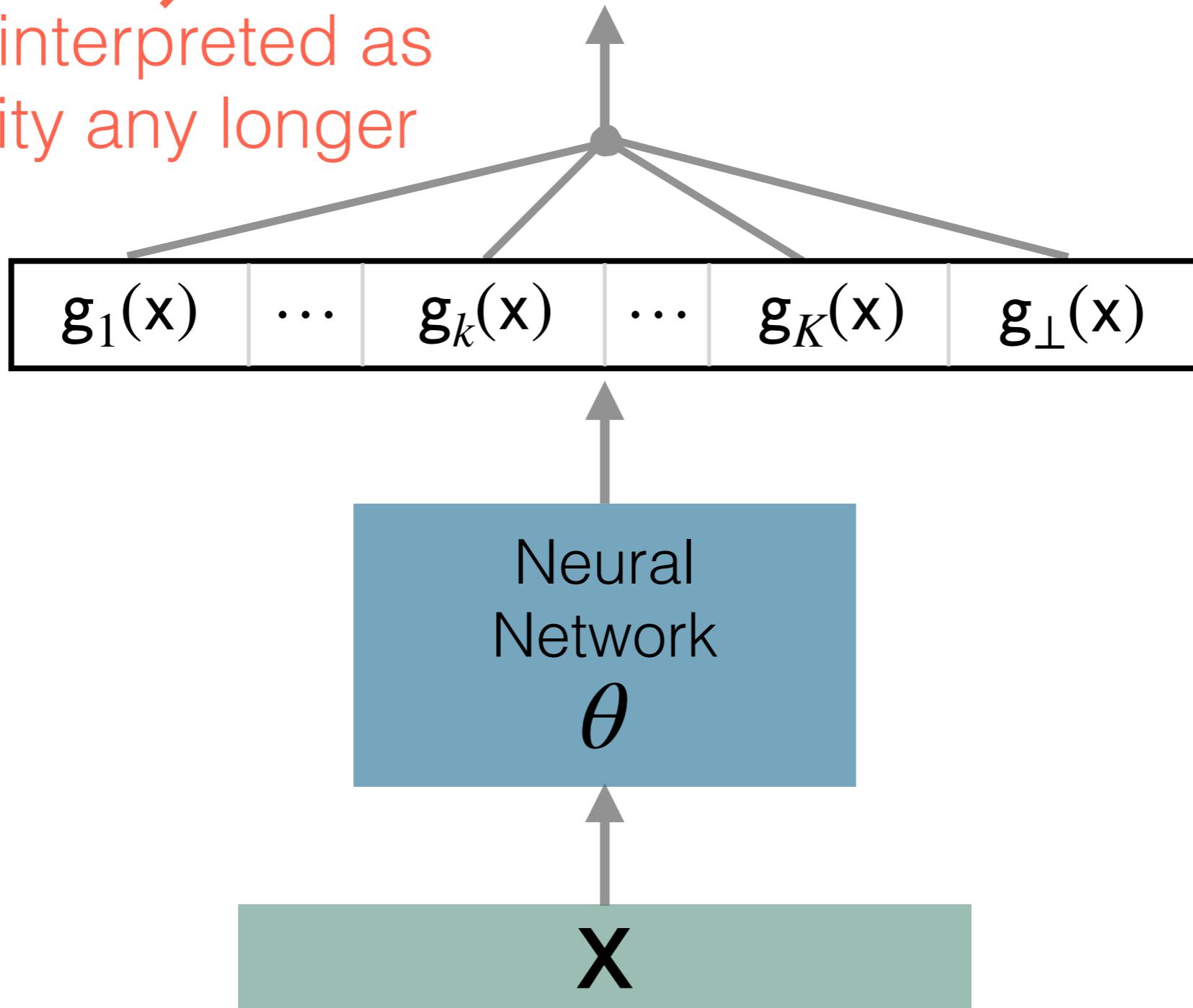


$$p_i(x) = \frac{\exp\{g_i(x)\}}{\sum_{k=1}^{K+1} \exp\{g_k(x)\}}$$

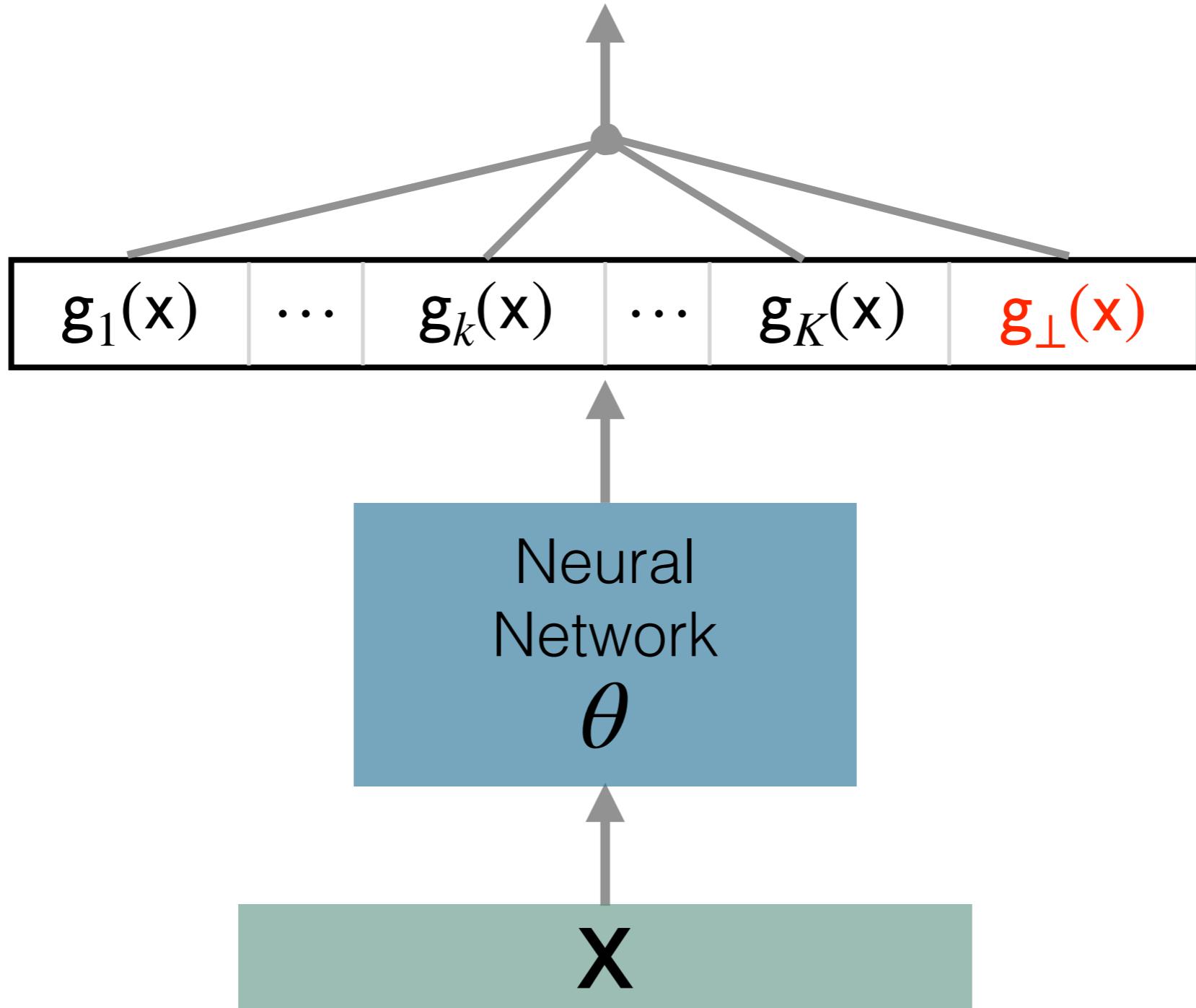


$$p_i(x) = \frac{\exp\{g_i(x)\}}{\sum_{k=1}^{K+1} \exp\{g_k(x)\}}$$

cannot be interpreted as  
a probability any longer



$$p_{\perp}(x) = \frac{\exp\{g_{\perp}(x)\}}{\sum_{k=1}^{K+1} \exp\{g_k(x)\}}$$



$$\ell(\theta; \mathcal{D}) =$$

$$\ell(\theta; \mathcal{D}) =$$

$$\ell(\theta; \mathcal{D}) = -\sum_n \left( \log p_{y_n}(x_n) + \mathbb{I}[y_n = m_n] \log p_{\perp}(x_n) \right)$$

$$\ell(\theta; \mathcal{D}) =$$

$$-\sum_n \left( \log p_{y_n}(x_n) + \mathbb{I}[y_n = m_n] \log p_{\perp}(x_n) \right)$$

classifier loss

rejector loss

only if expert is correct

$$\ell(\theta; \mathcal{D}) =$$

$$-\sum_n \left( \log p_{y_n}(x_n) + \mathbb{I}[y_n = m_n] \log p_{\perp}(x_n) \right)$$

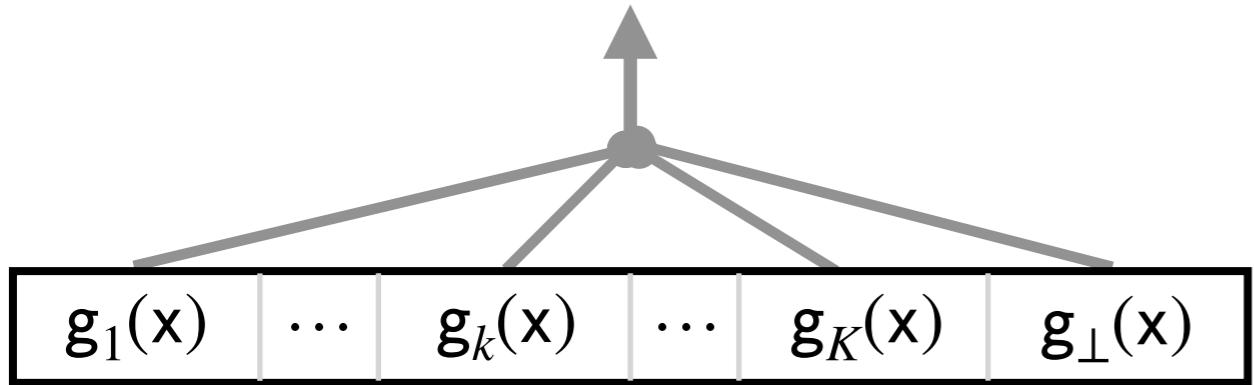
classifier loss

rejector loss

only if expert is correct

**Consistency:** The minimizers (w.r.t.  $g$ ) correspond to the Bayes optimal classifier and rejector.

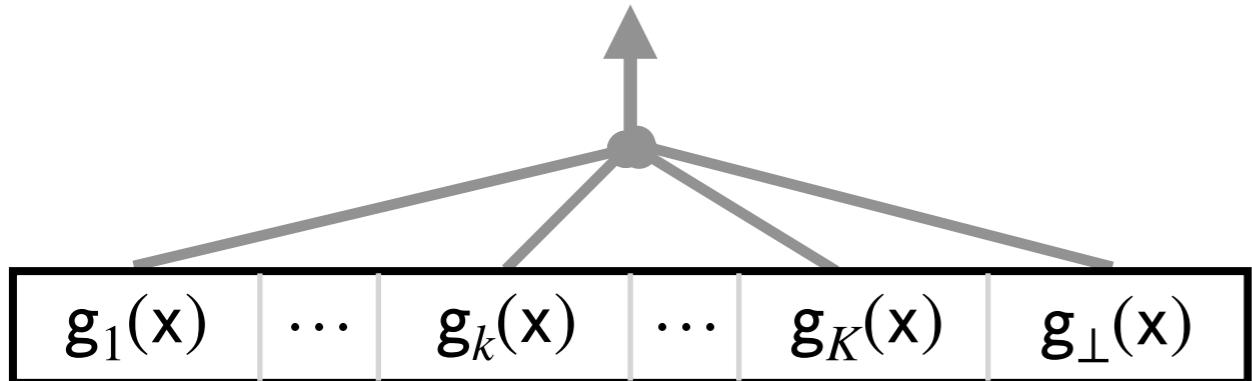
$$p_i(x) = \frac{\exp\{g_i(x)\}}{\sum_{k=1}^{K+1} \exp\{g_k(x)\}}$$



# Classifier

$$h^*(x) = \operatorname{argmax}_{y \in [1, K]} g_y^*(x)$$

$$p_i(x) = \frac{\exp\{g_i(x)\}}{\sum_{k=1}^{K+1} \exp\{g_k(x)\}}$$



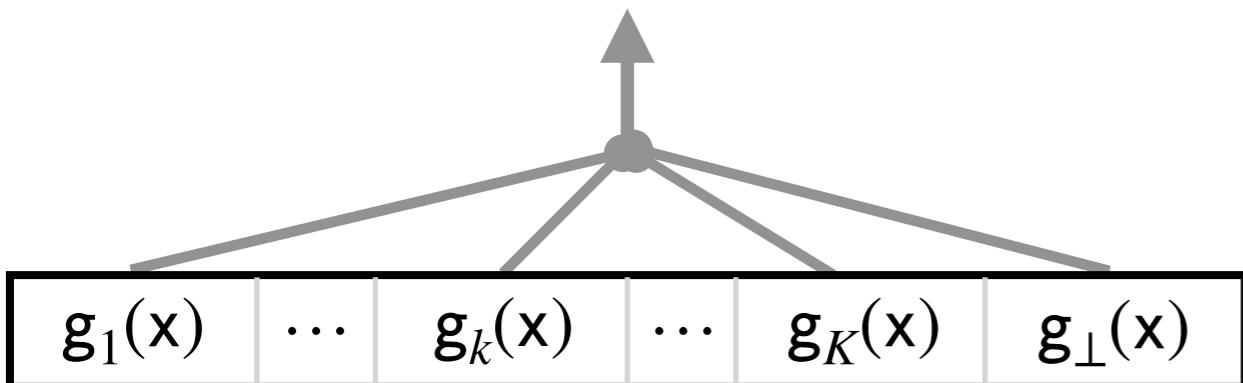
# Classifier

$$h^*(x) = \operatorname{argmax}_{y \in [1, K]} g_y^*(x)$$

# Rejector

$$r^*(x) = \mathbb{I} \left[ g_{\perp}^*(x) \geq \max_{y \in [1, K]} g_y^*(x) \right]$$

$$p_i(x) = \frac{\exp\{g_i(x)\}}{\sum_{k=1}^{K+1} \exp\{g_k(x)\}}$$



Our work:  
Is this learning-to-defer  
system calibrated?



# Human-computer collaboration for skin cancer recognition

Philipp Tschandl<sup>1,17</sup>, Christoph Rinner<sup>1,17</sup>, Zoe Apalla<sup>3</sup>, Giuseppe Argenziano<sup>1,4</sup>, Noel Codella<sup>5</sup>, Allan Halpern<sup>6</sup>, Monika Janda<sup>7</sup>, Aimilios Lallas<sup>3</sup>, Caterina Longo<sup>8,9</sup>, Josep Malvehy<sup>10,11</sup>, John Paoli<sup>12,13</sup>, Susana Puig<sup>10,11</sup>, Cliff Rosendahl<sup>14</sup>, H. Peter Soyer<sup>15</sup>, Iris Zalaudek<sup>16</sup> and Harald Kittler<sup>1</sup>✉

The rapid increase in telemedicine coupled with recent advances in diagnostic artificial intelligence (AI) create the imperative to consider the opportunities and risks of inserting AI-based support into new paradigms of care. Here we

competitive view of AI is evolving based on studies suggesting that a more promising approach is human–AI cooperation<sup>10–15</sup>. The role of human–computer collaboration in health-care delivery, the appropriate settings in which it can be applied and its impact on the

“The least experienced [physicians] tended to accept AI-based support that contradicted their initial diagnosis even if they were confident.”

# Calibration: Is the system a good forecaster?

# Calibration: Is the system a good forecaster?

$$P(y | x)$$

Does the classifier correctly estimate the underlying class probabilities?

# Calibration: Is the system a good forecaster?

$$P(y | x)$$

Does the classifier correctly estimate the underlying class probabilities?

$$P(m = y | x)$$

Does the rejector correctly estimate the expert's chance of being correct?

# Calibration: Is the system a good forecaster?

$$P(y | x)$$



Does the classifier correctly estimate the underlying class probabilities?

$$P(m = y | x)$$

Does the rejector correctly estimate the expert's chance of being correct?

# Calibration: Is the system a good forecaster?

$$P(y | x)$$



Does the classifier correctly estimate the underlying class probabilities?

$$P(m = y | x)$$



Does the rejector correctly estimate the expert's chance of being correct?

# Calibration: Is the system a good forecaster?

$$P(y | x)$$



Does the classifier correctly estimate the underlying class probabilities?

$$P(m = y | x)$$



Does the rejector correctly estimate the expert's chance of being correct?

Recall...

$$p_{\perp}(x) = \frac{\exp\{g_{\perp}(x)\}}{\sum_{k=1}^{K+1} \exp\{g_k(x)\}}$$

Recall...

$$p_{\perp}(x) = \frac{\exp\{g_{\perp}(x)\}}{\sum_{k=1}^{K+1} \exp\{g_k(x)\}}$$

From Mozannar's & Sontag's Theorem #1...

$$\frac{\mathbb{P}(m = y | x)}{1 + \mathbb{P}(m = y | x)} = p_{\perp}^*(x)$$

Recall...

$$p_{\perp}(x) = \frac{\exp\{g_{\perp}(x)\}}{\sum_{k=1}^{K+1} \exp\{g_k(x)\}}$$

From Mozannar's & Sontag's Theorem #1...

$$\frac{\mathbb{P}(m = y | x)}{1 + \mathbb{P}(m = y | x)} = p_{\perp}^*(x)$$

Rearranging...

$$\mathbb{P}(m = y | x) = \frac{p_{\perp}^*(x)}{1 - p_{\perp}^*(x)}$$

Recall...

$$p_{\perp}(x) = \frac{\exp\{g_{\perp}(x)\}}{\sum_{k=1}^{K+1} \exp\{g_k(x)\}}$$

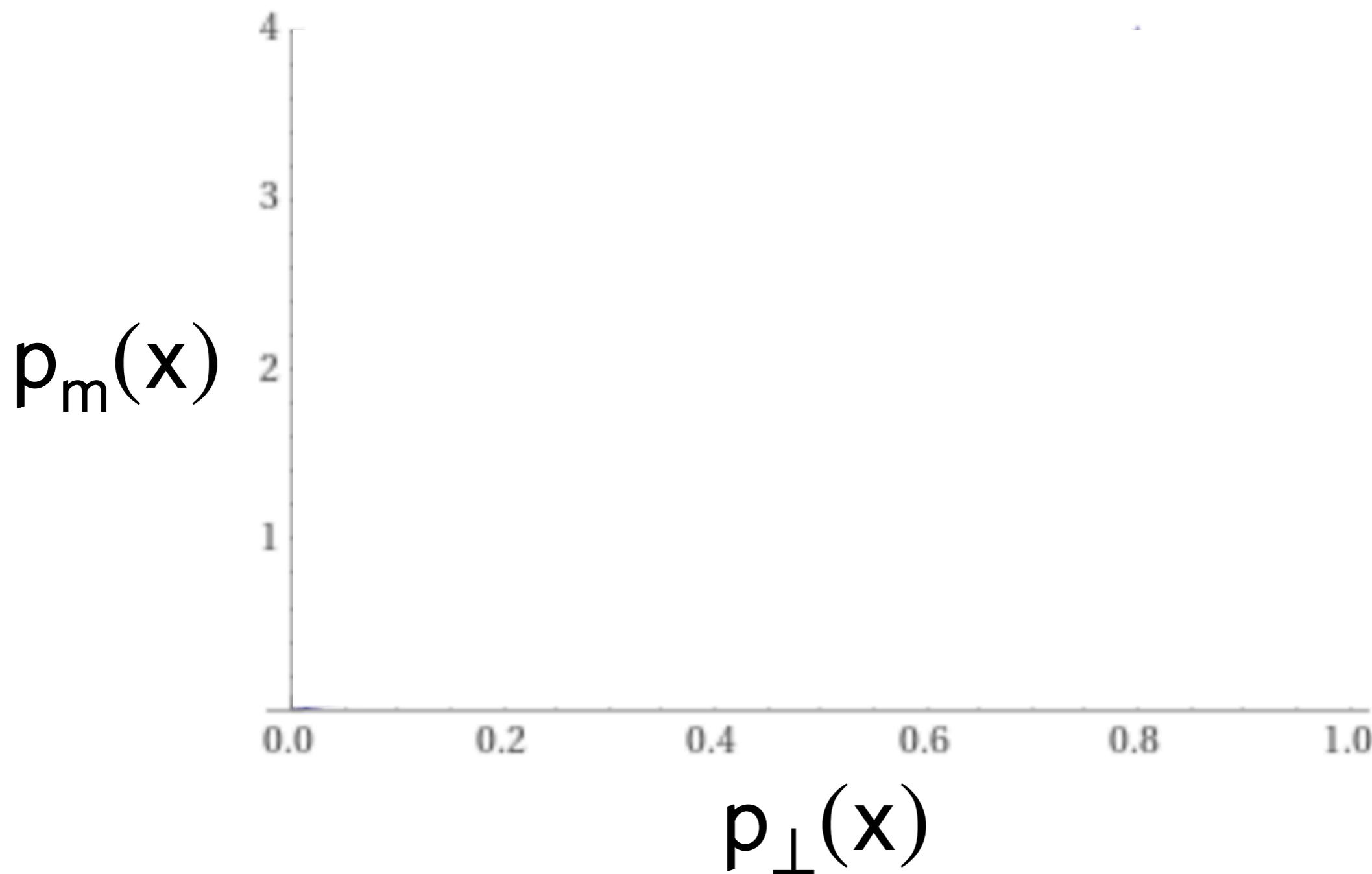
From Mozannar's & Sontag's Theorem #1...

$$\frac{\mathbb{P}(m = y | x)}{1 + \mathbb{P}(m = y | x)} = p_{\perp}^*(x)$$

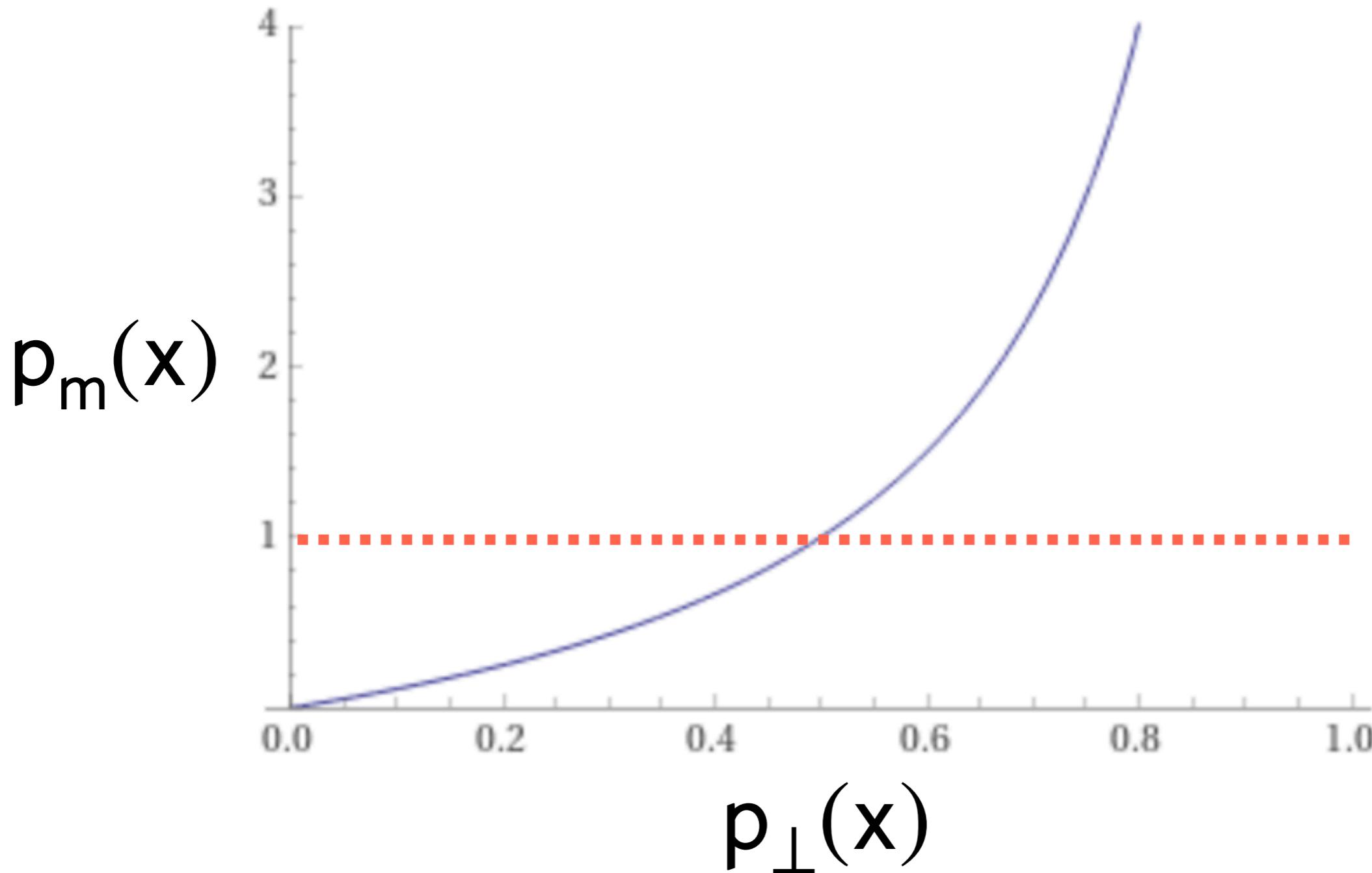
Rearranging...

$$\mathbb{P}(m = y | x) = \frac{p_{\perp}^*(x)}{1 - p_{\perp}^*(x)} \triangleq p_m(x)$$

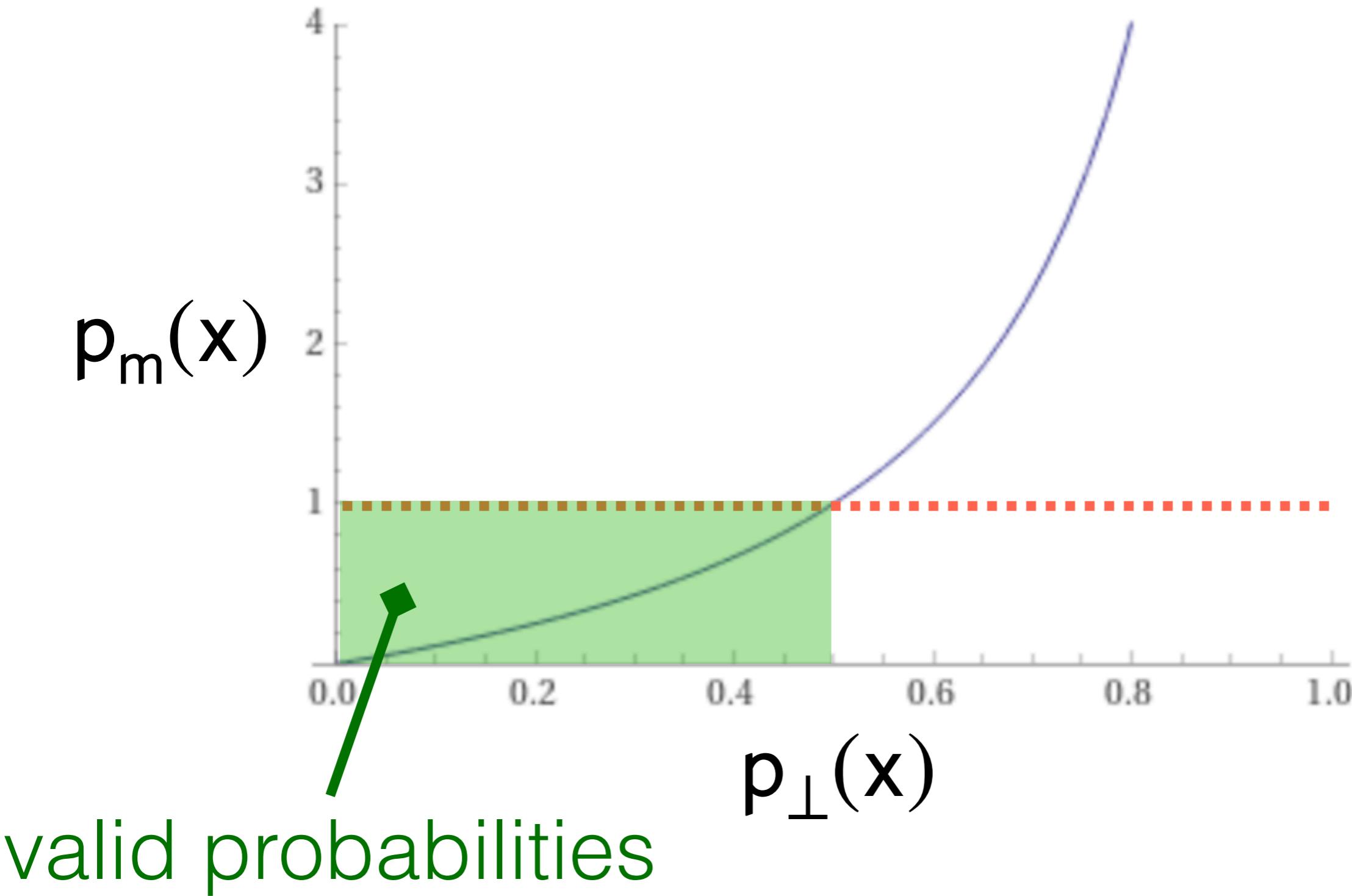
# Estimating $P(m = y | x)$



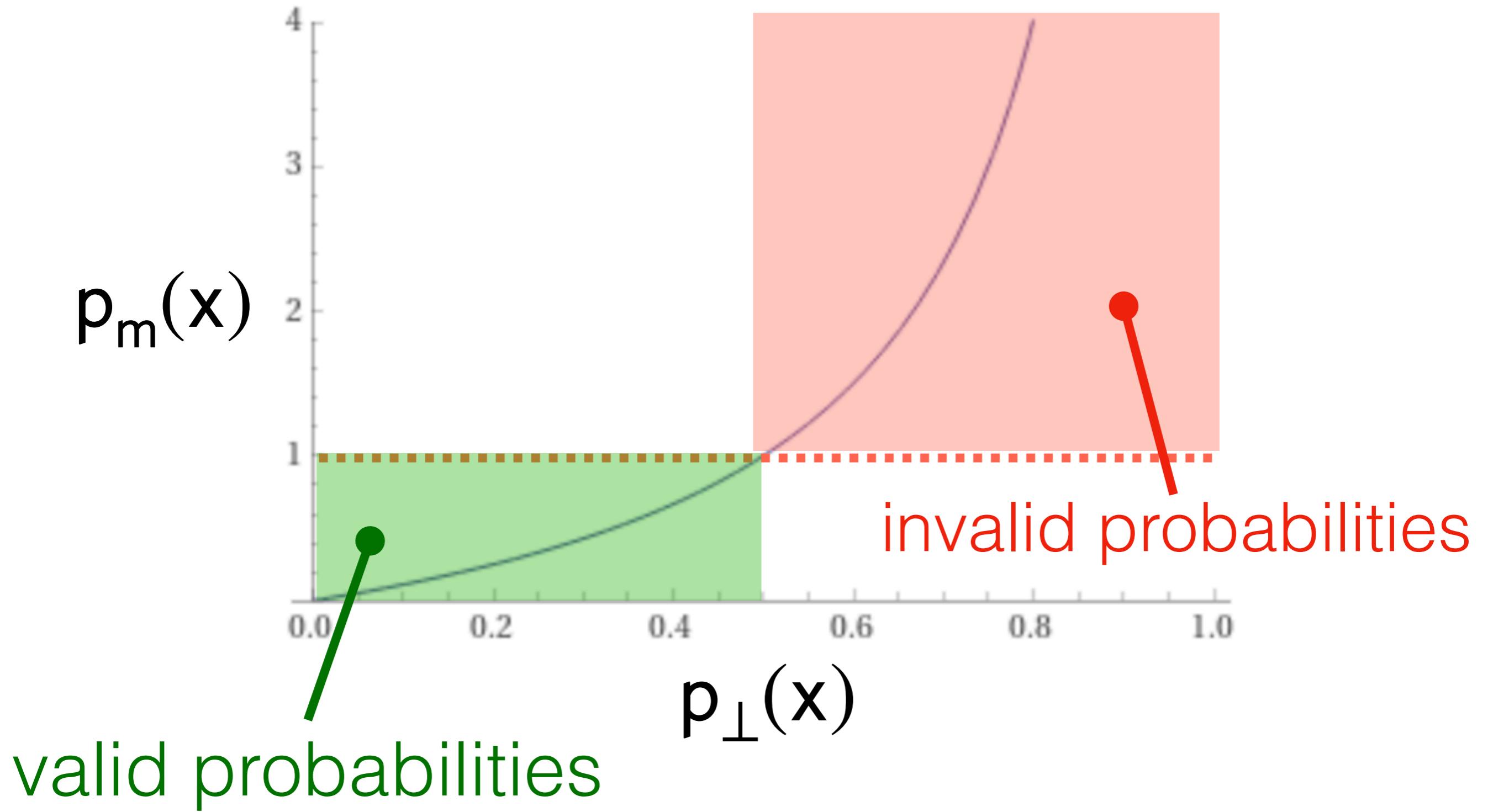
# Estimating $P(m = y | x)$



# Estimating $P(m = y | x)$



# Estimating $P(m = y | x)$

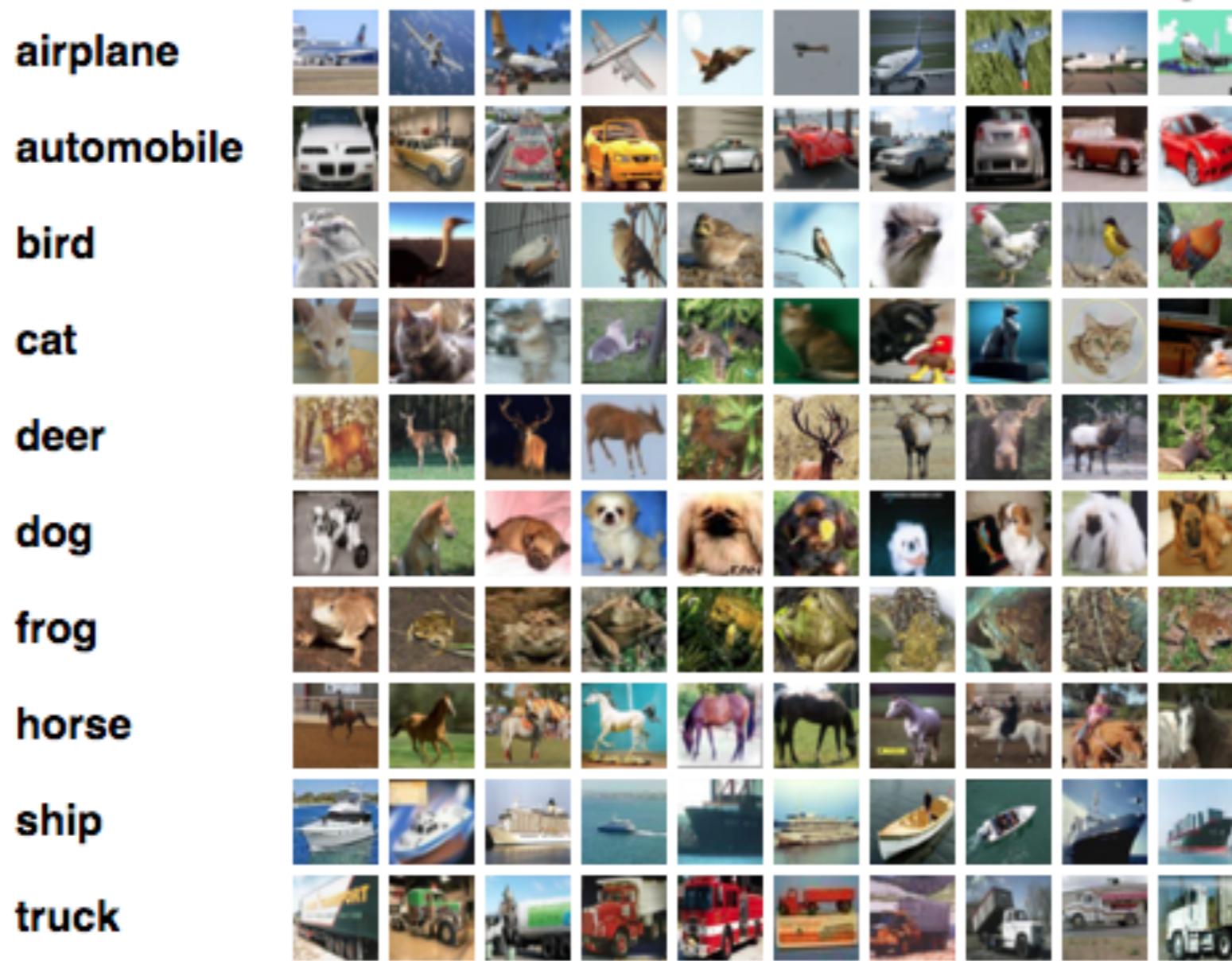


# What happens in practice?

CIFAR-10: Image Classification

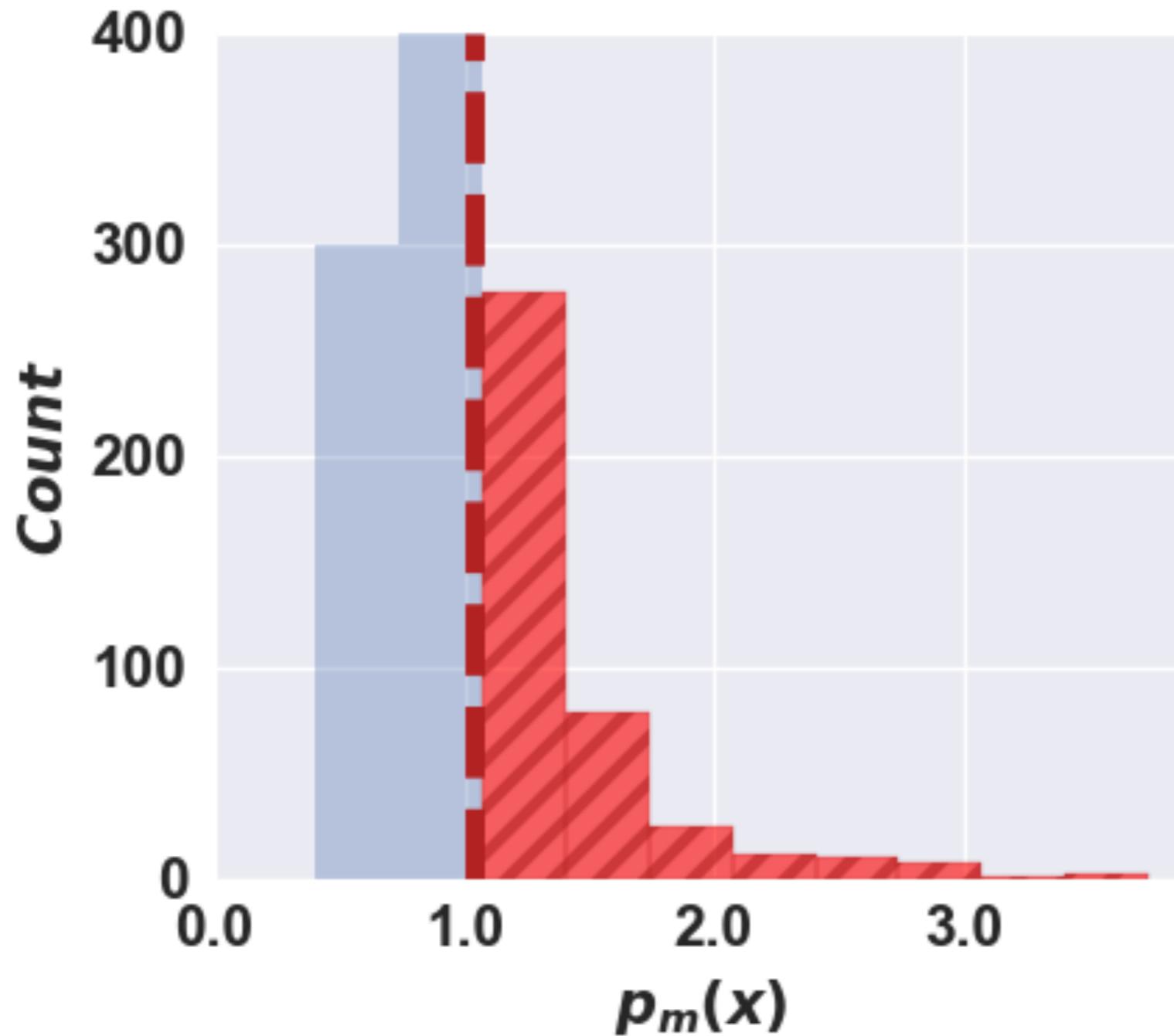
# What happens in practice?

## CIFAR-10: Image Classification



# What happens in practice?

CIFAR-10: Image Classification

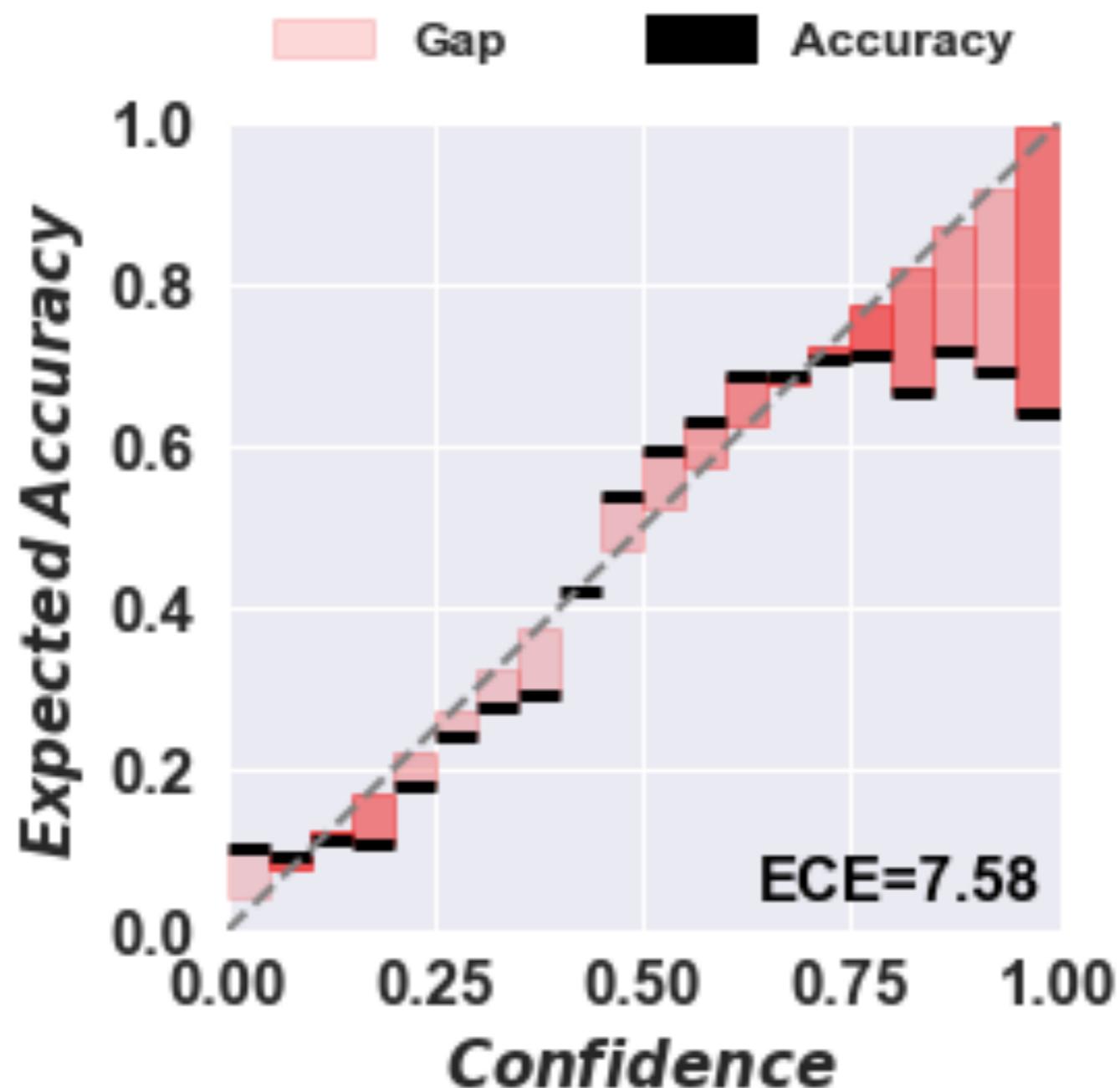


Truncation:  $p_m(x) \in (0,1]$

CIFAR-10: Image Classification

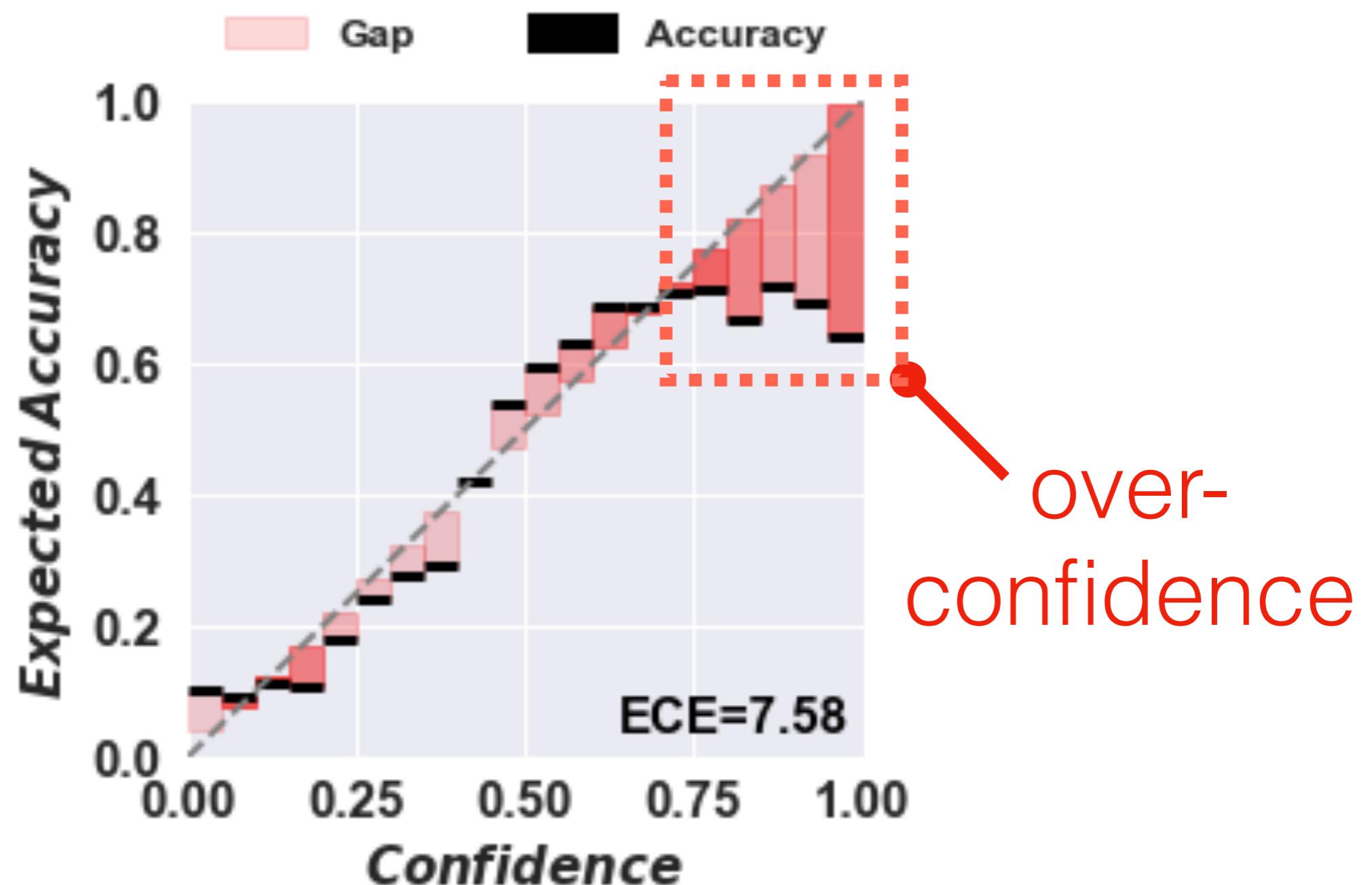
Truncation:  $p_m(x) \in (0,1]$

CIFAR-10: Image Classification



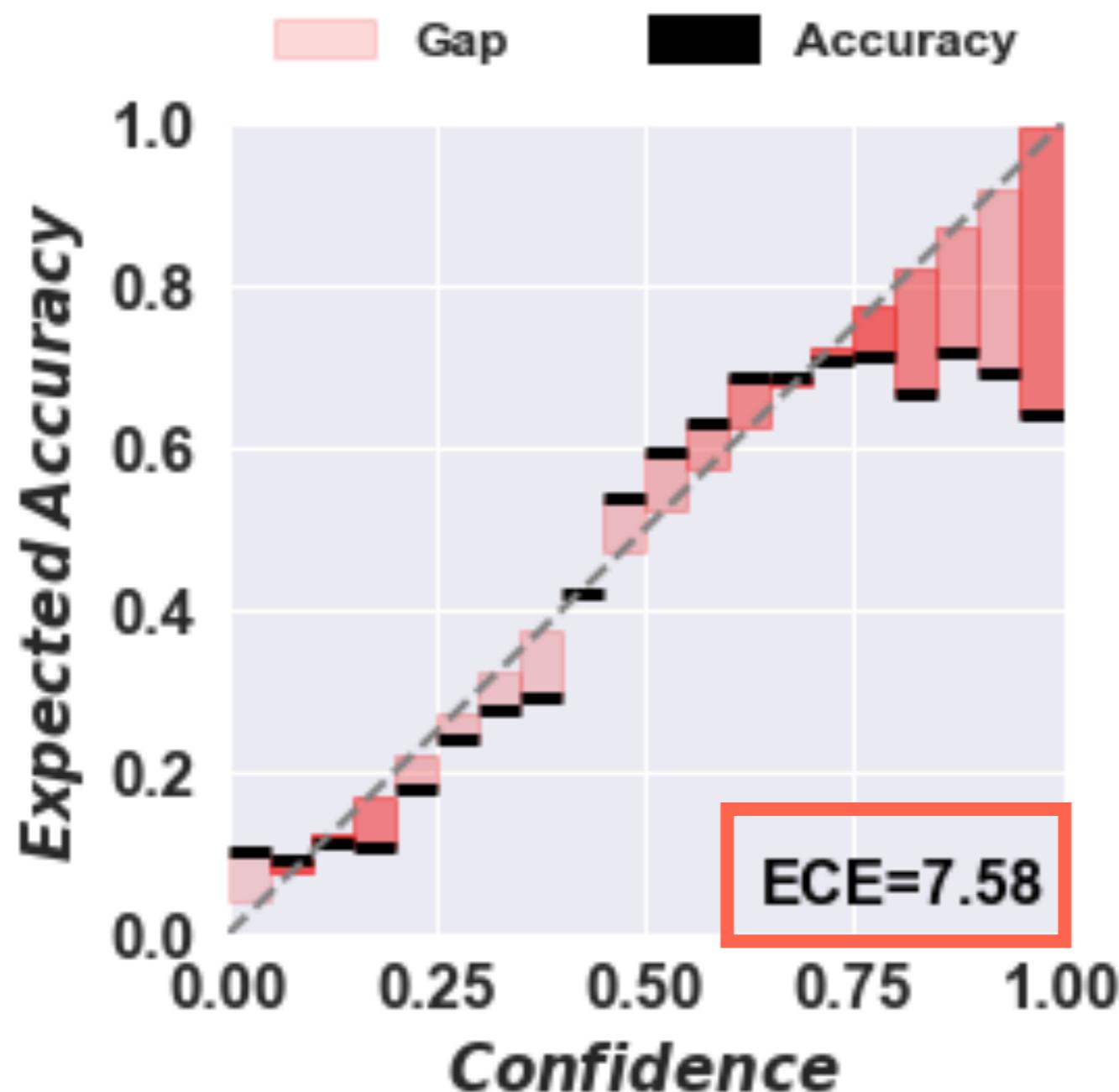
Truncation:  $p_m(x) \in (0,1]$

CIFAR-10: Image Classification



Truncation:  $p_m(x) \in (0,1]$

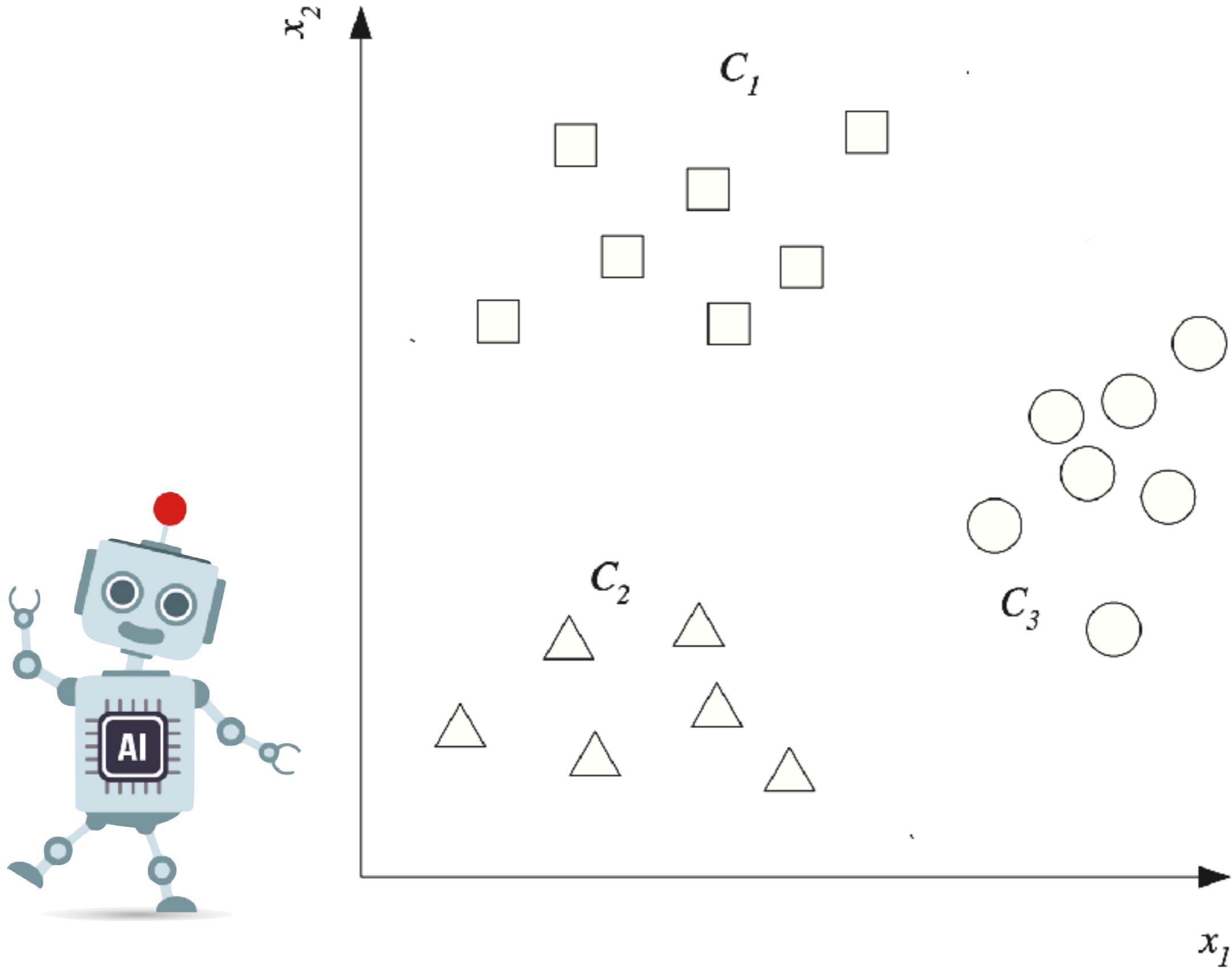
CIFAR-10: Image Classification



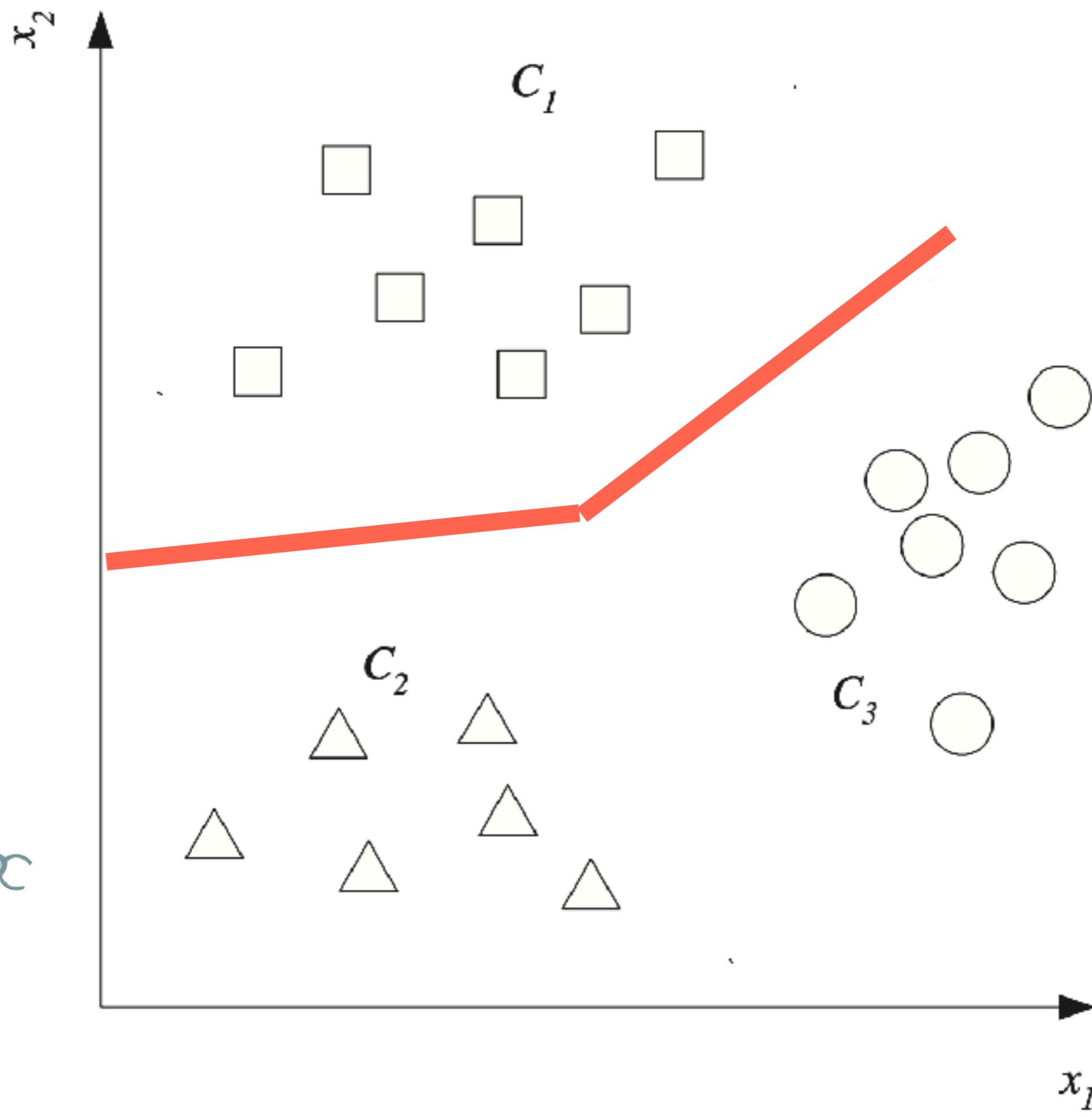
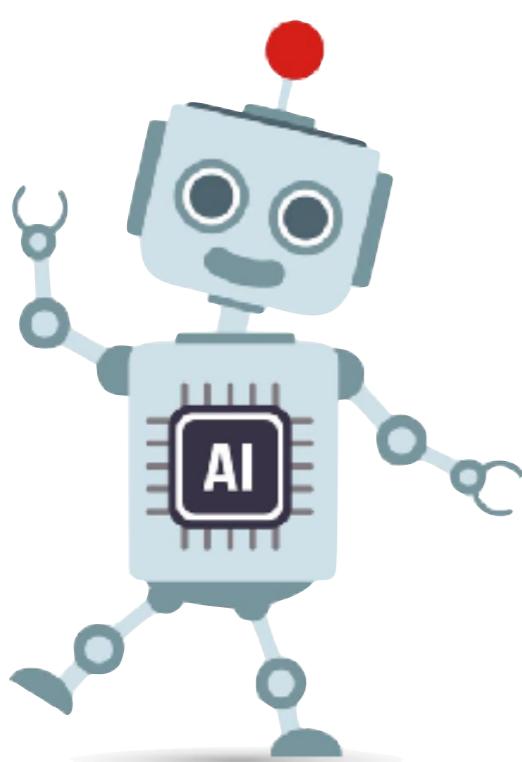
$ECE =$   
expected  
calibration  
error

# Solution: Learning to Defer with a One-vs-All Classifier

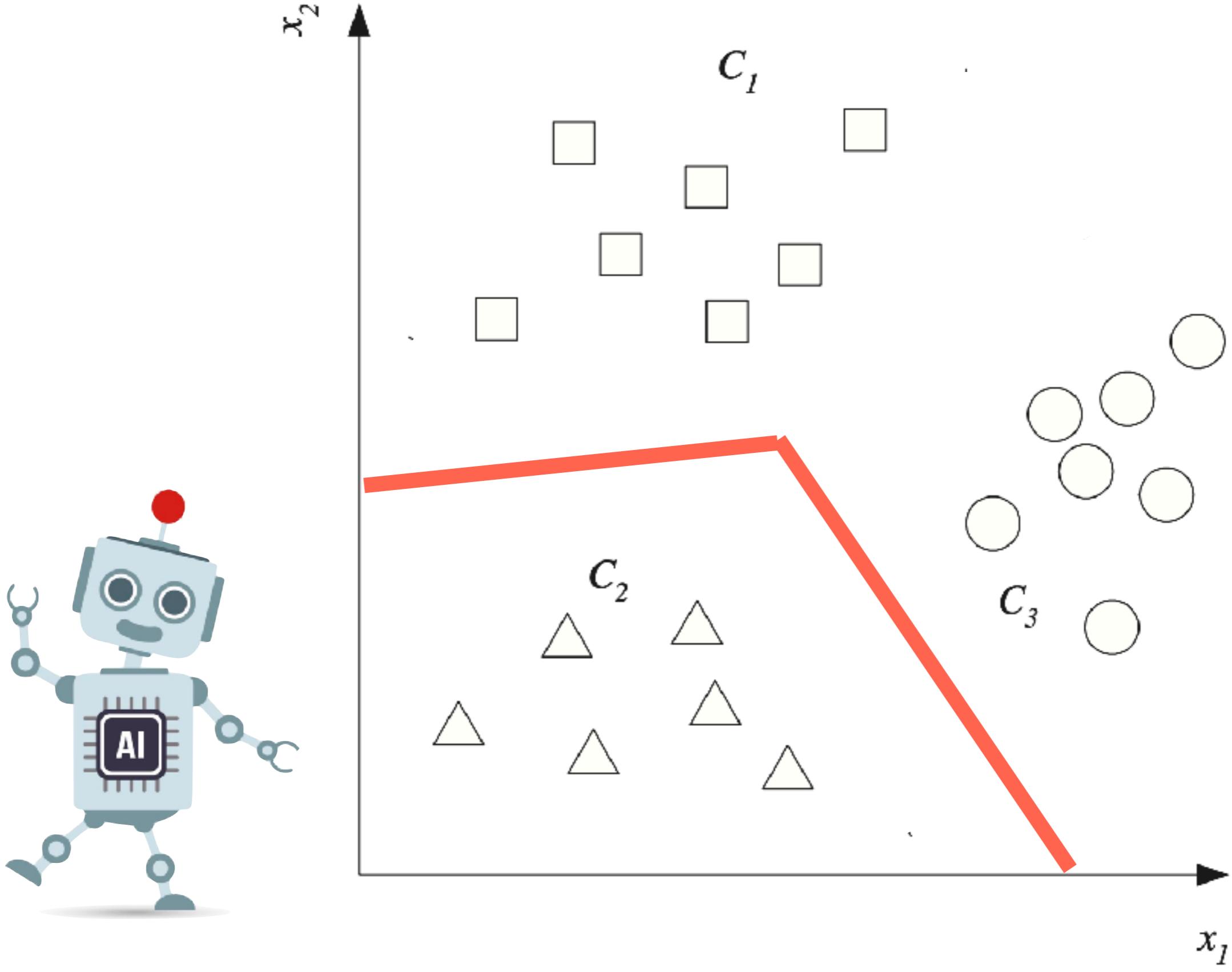
# Solution: Learning to Defer with a One-vs-All Classifier



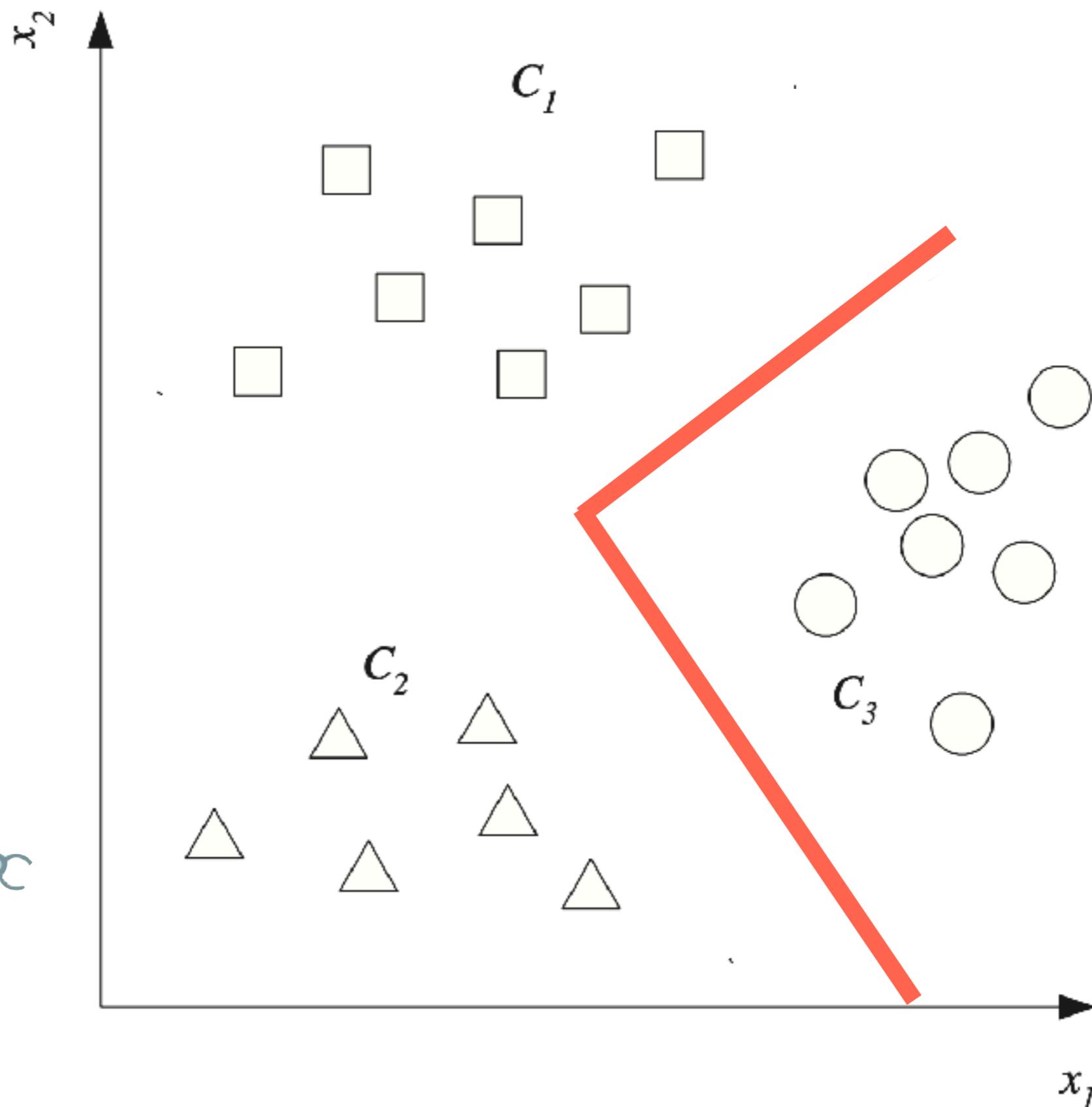
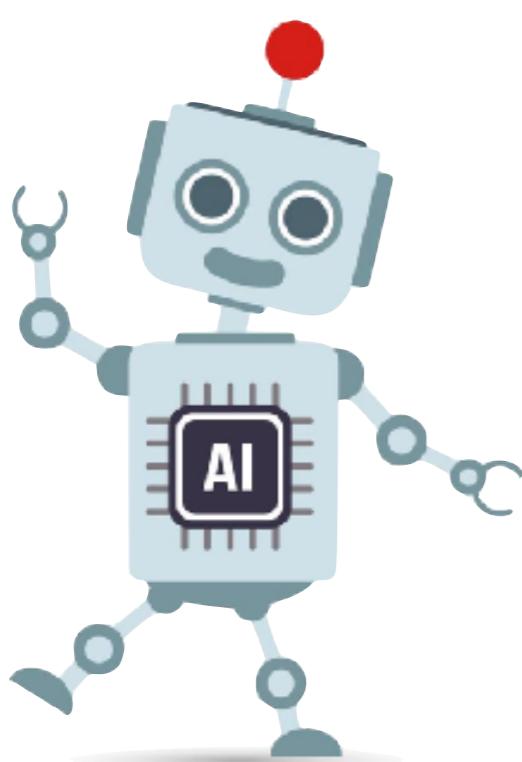
One-vs-All Classifier



One-vs-All Classifier



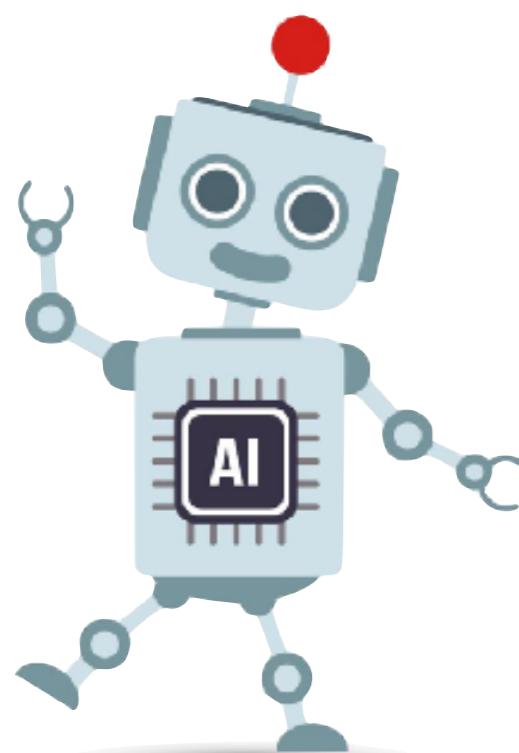
One-vs-All Classifier



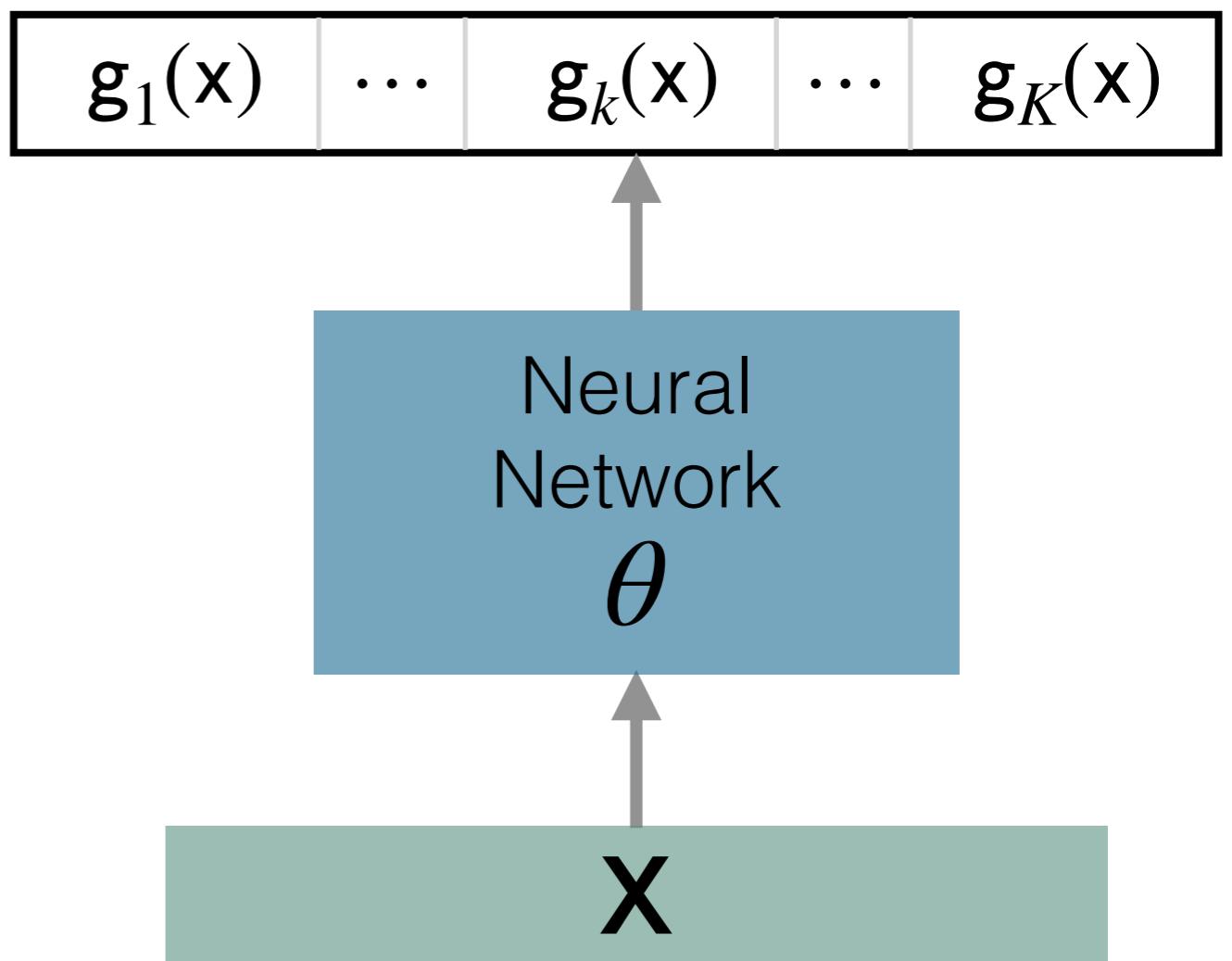
One-vs-All Classifier

100

**Goal:** For input features  $\mathbf{x}$ ,  
predict membership in one  
of  $K$  classes:  $C_1, \dots, C_K$

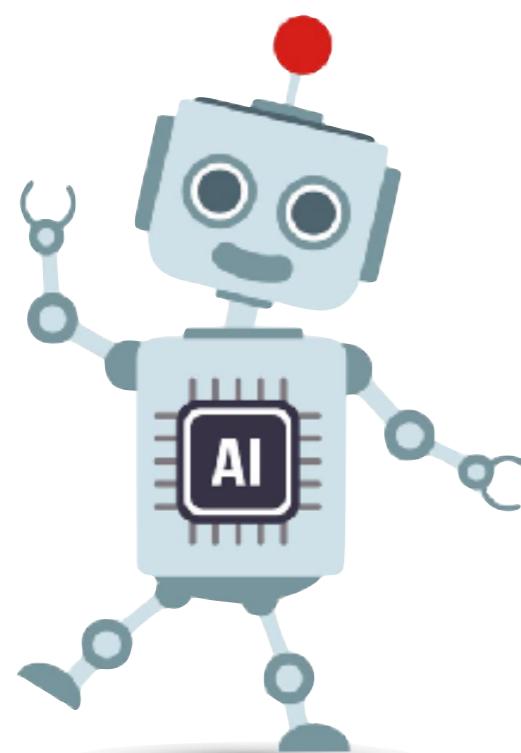


One-vs-All Classifier

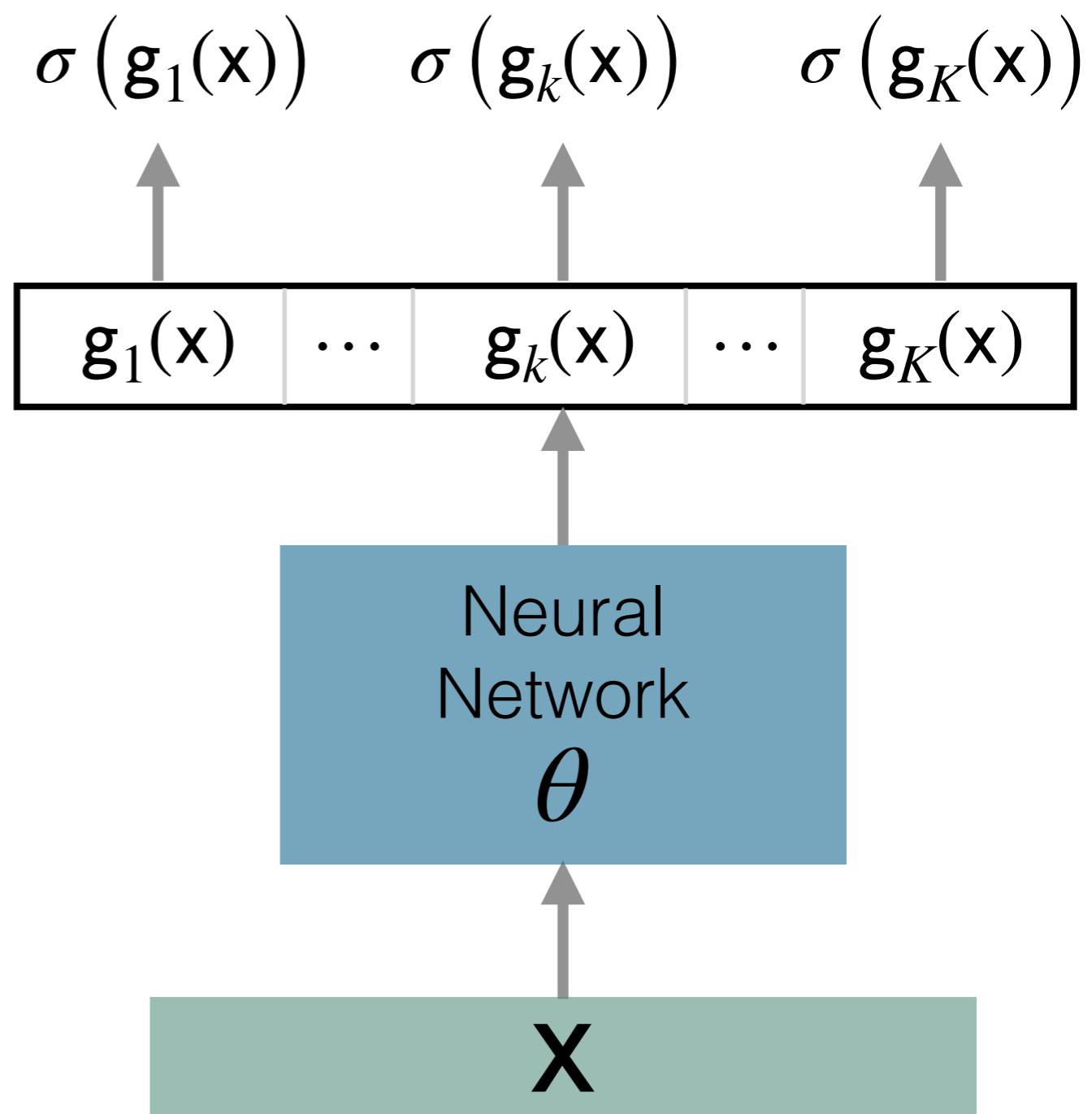


Goal: For input features  $\mathbf{x}$ ,  
predict membership in one  
of  $K$  classes:  $C_1, \dots, C_K$

$$\sigma(z) = \frac{1}{1 + \exp\{-z\}}$$

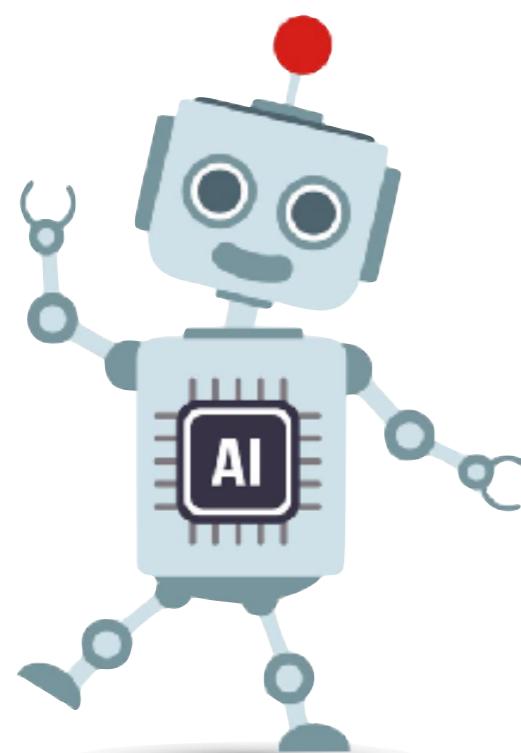


One-vs-All Classifier

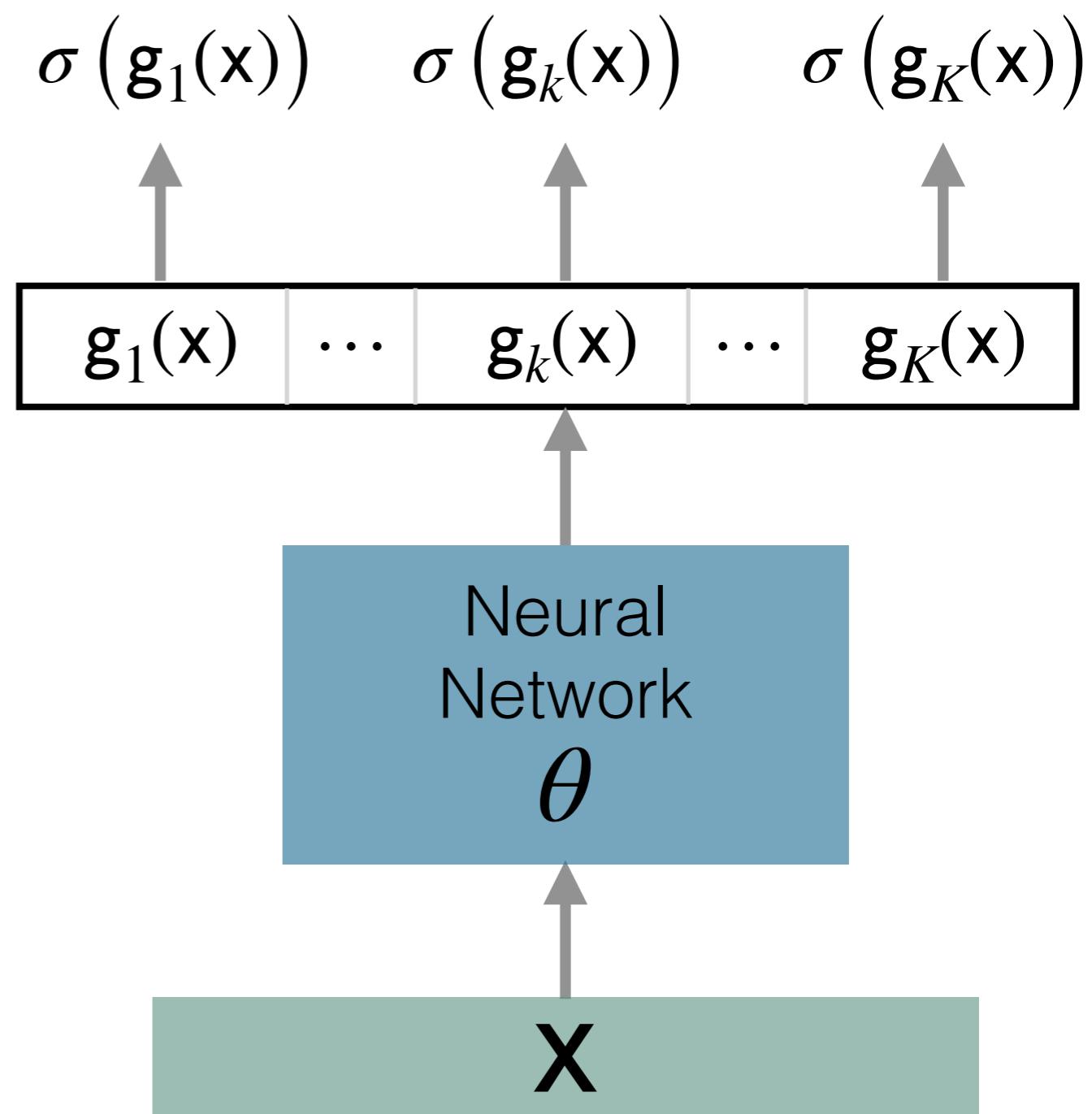


Goal: For input features  $\mathbf{x}$ ,  
predict membership in one  
of  $K$  classes:  $C_1, \dots, C_K$

$$P(C_i | \mathbf{x}) = \sigma(g_i(\mathbf{x}))$$

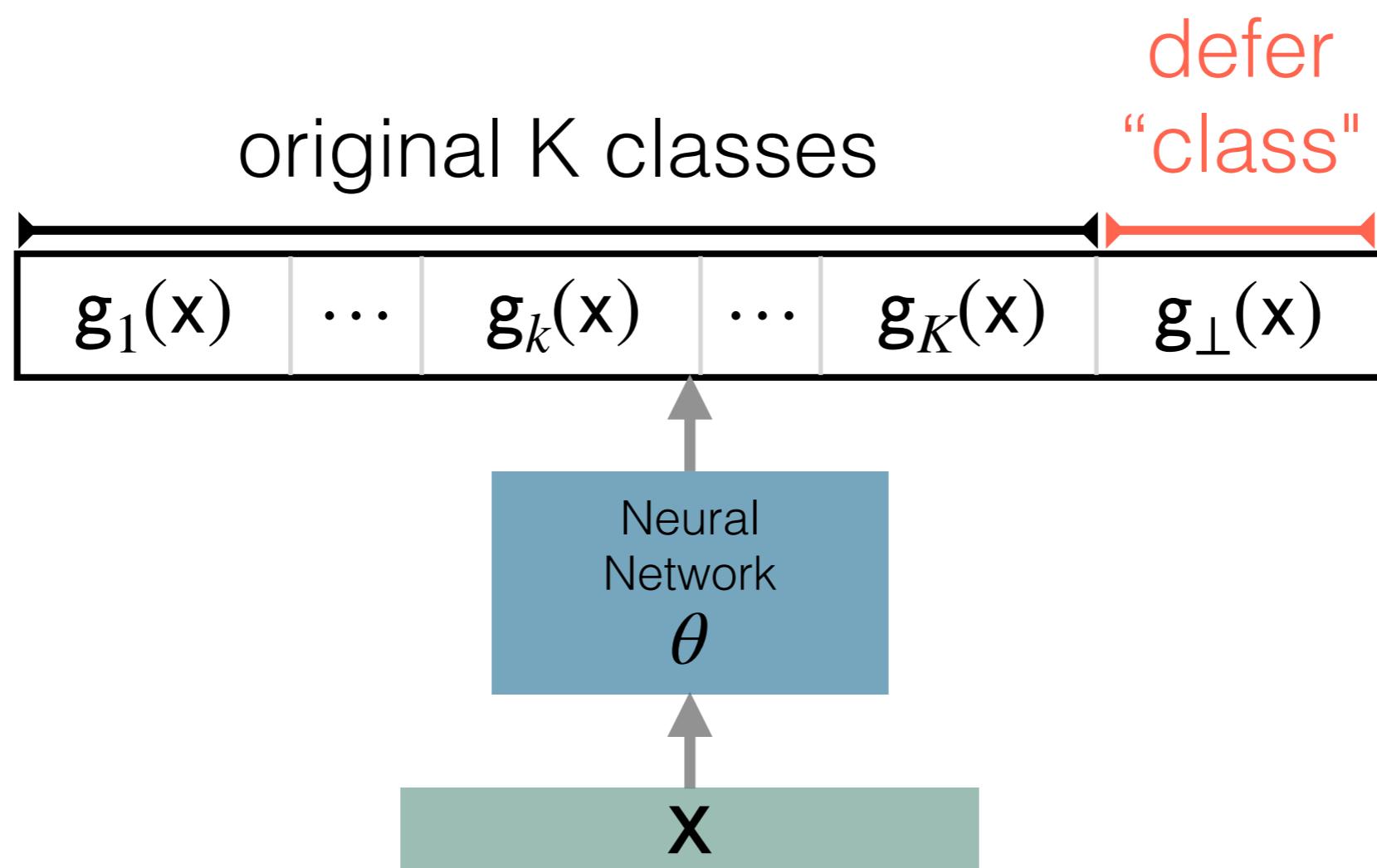


One-vs-All Classifier

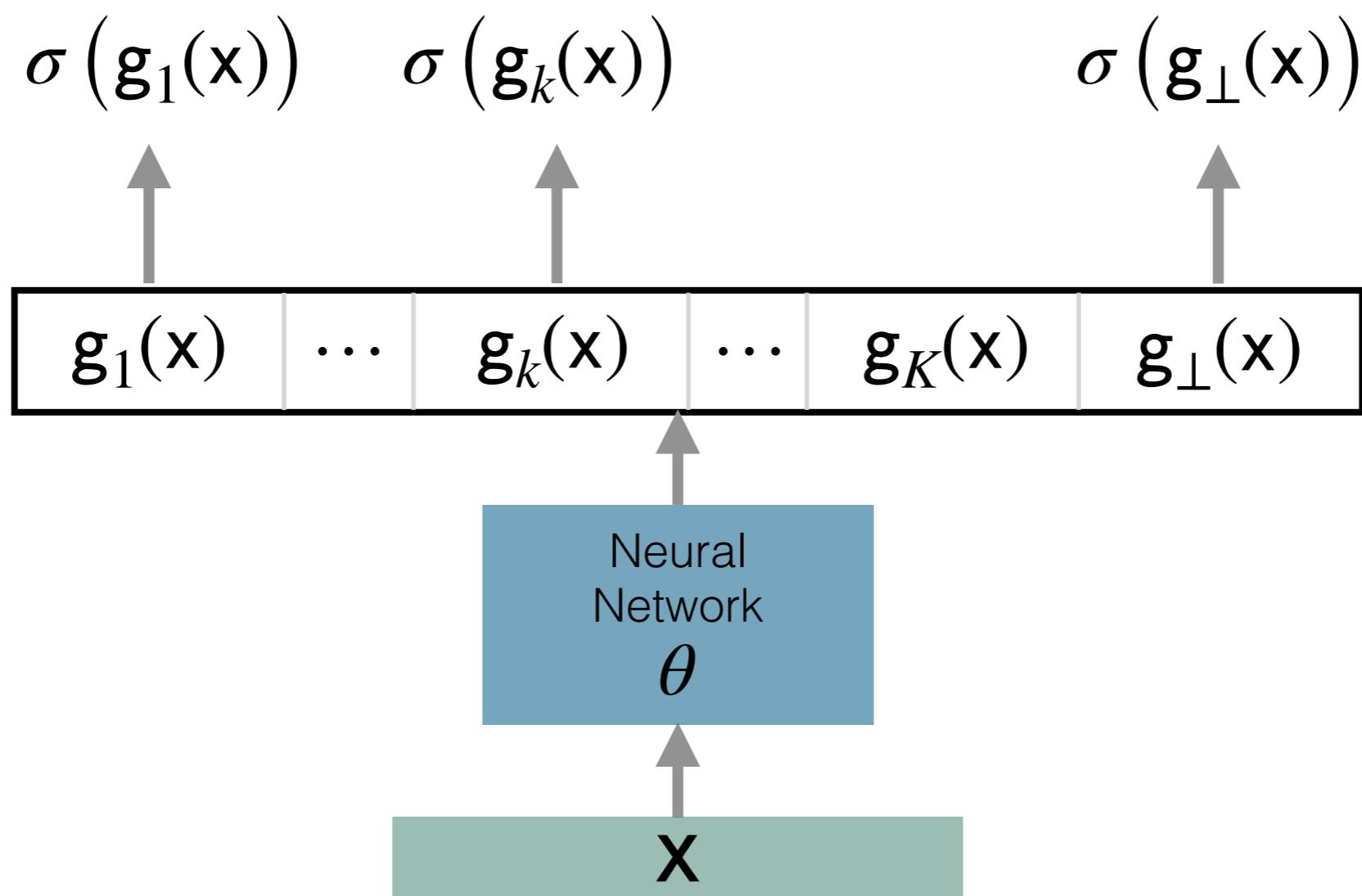


# Learning-to-Defer: One-vs-All Parameterization

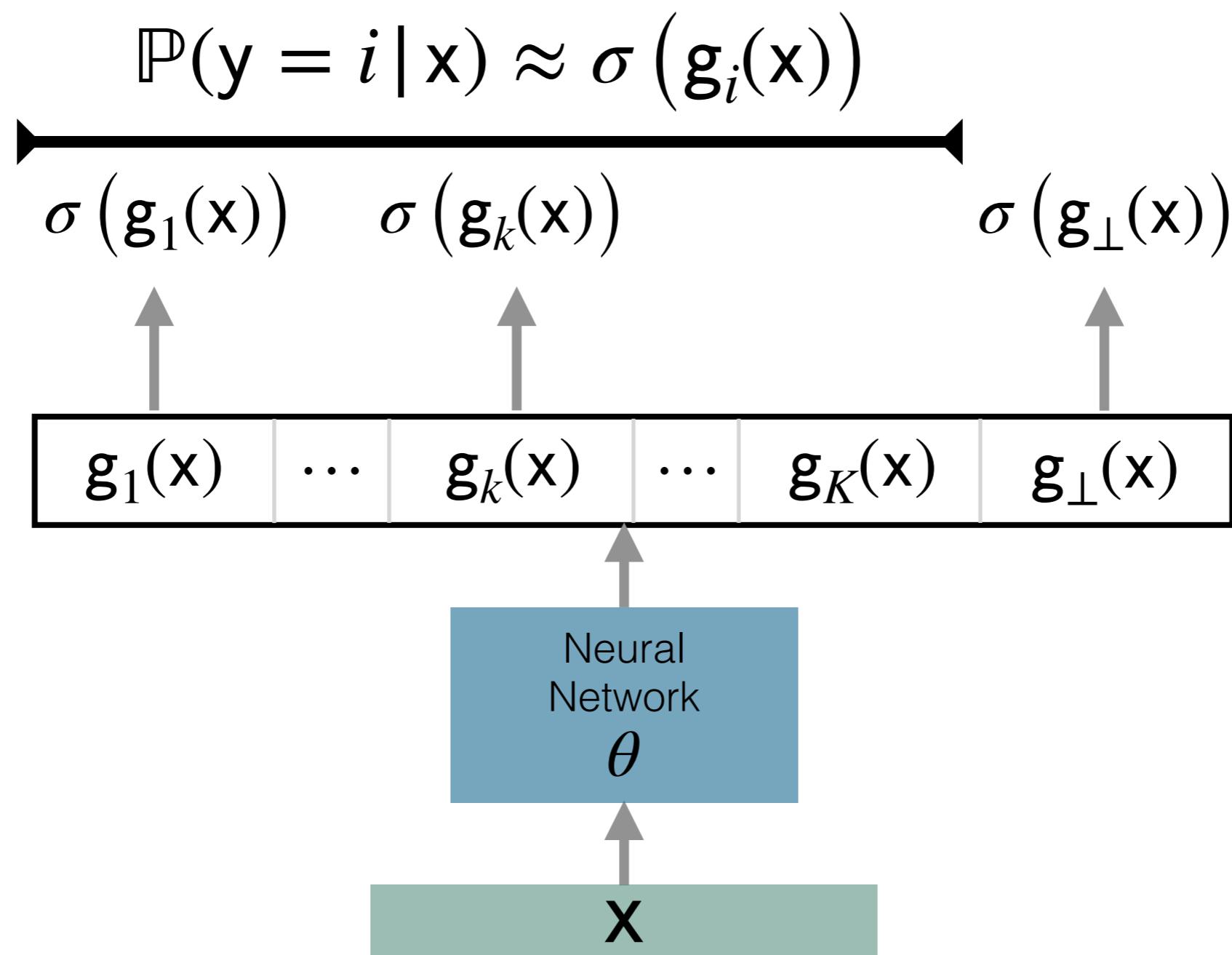
# Learning-to-Defer: One-vs-All Parameterization



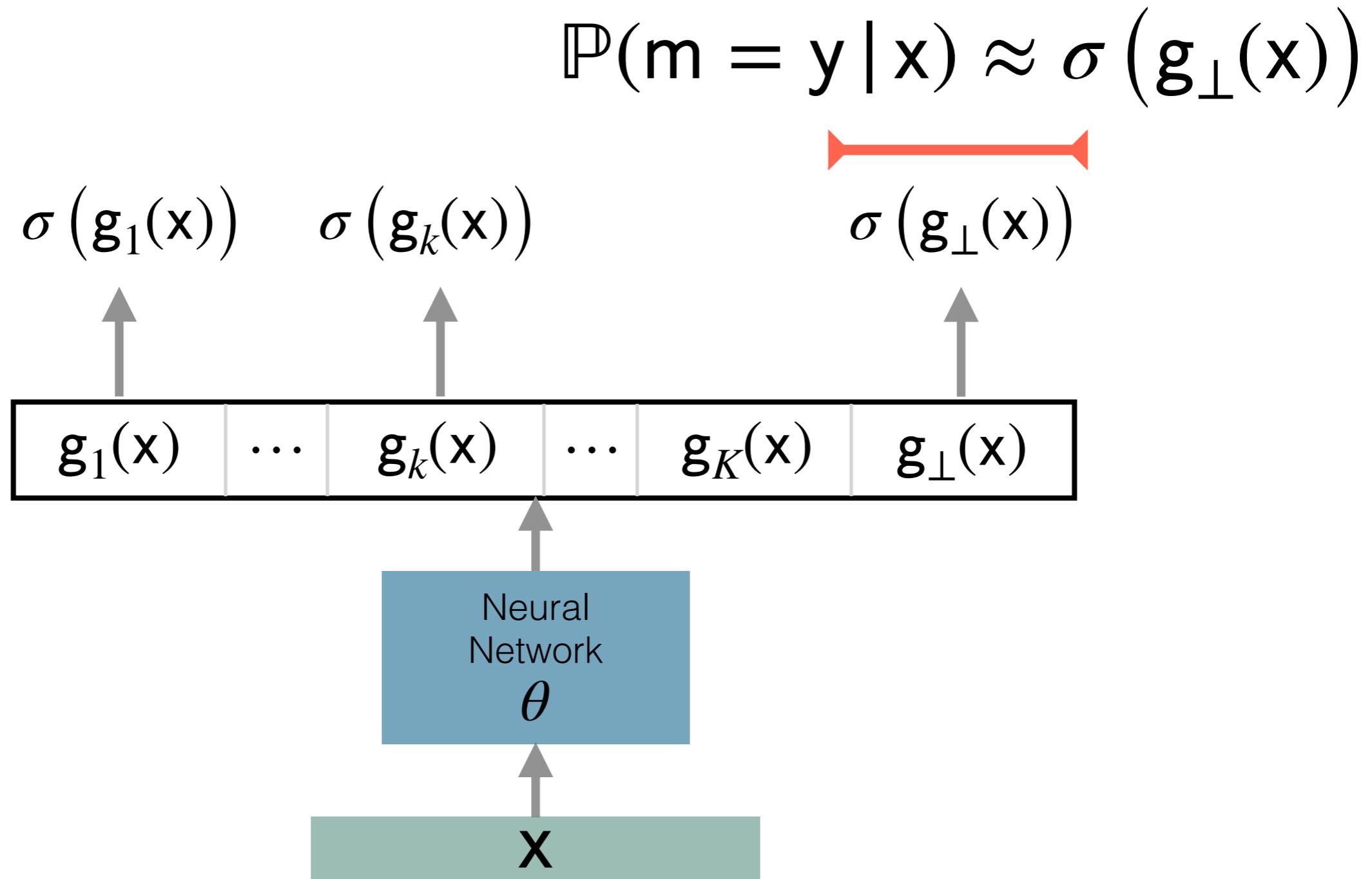
# Learning-to-Defer: One-vs-All Parameterization



# Learning-to-Defer: One-vs-All Parameterization



# Learning-to-Defer: One-vs-All Parameterization



# Learning-to-Defer: One-vs-All Parameterization

**Theorem 4.1 + Corollary 4.2** [Verma & Nalisnick, 2022]:  
The one-vs-all loss is a *consistent* surrogate for  
the 0-1 learning-to-defer loss.

# Learning-to-Defer: One-vs-All Parameterization

**Theorem 4.1 + Corollary 4.2** [Verma & Nalisnick, 2022]:  
The one-vs-all loss is a *consistent* surrogate for  
the 0-1 learning-to-defer loss.

**Proof:** Via method of *error correcting output codes* [Dietterich & Bakiri, 1995], by applying  
Ramaswamy et al.'s [2018] Equation 1.

# Learning-to-Defer: One-vs-All Parameterization

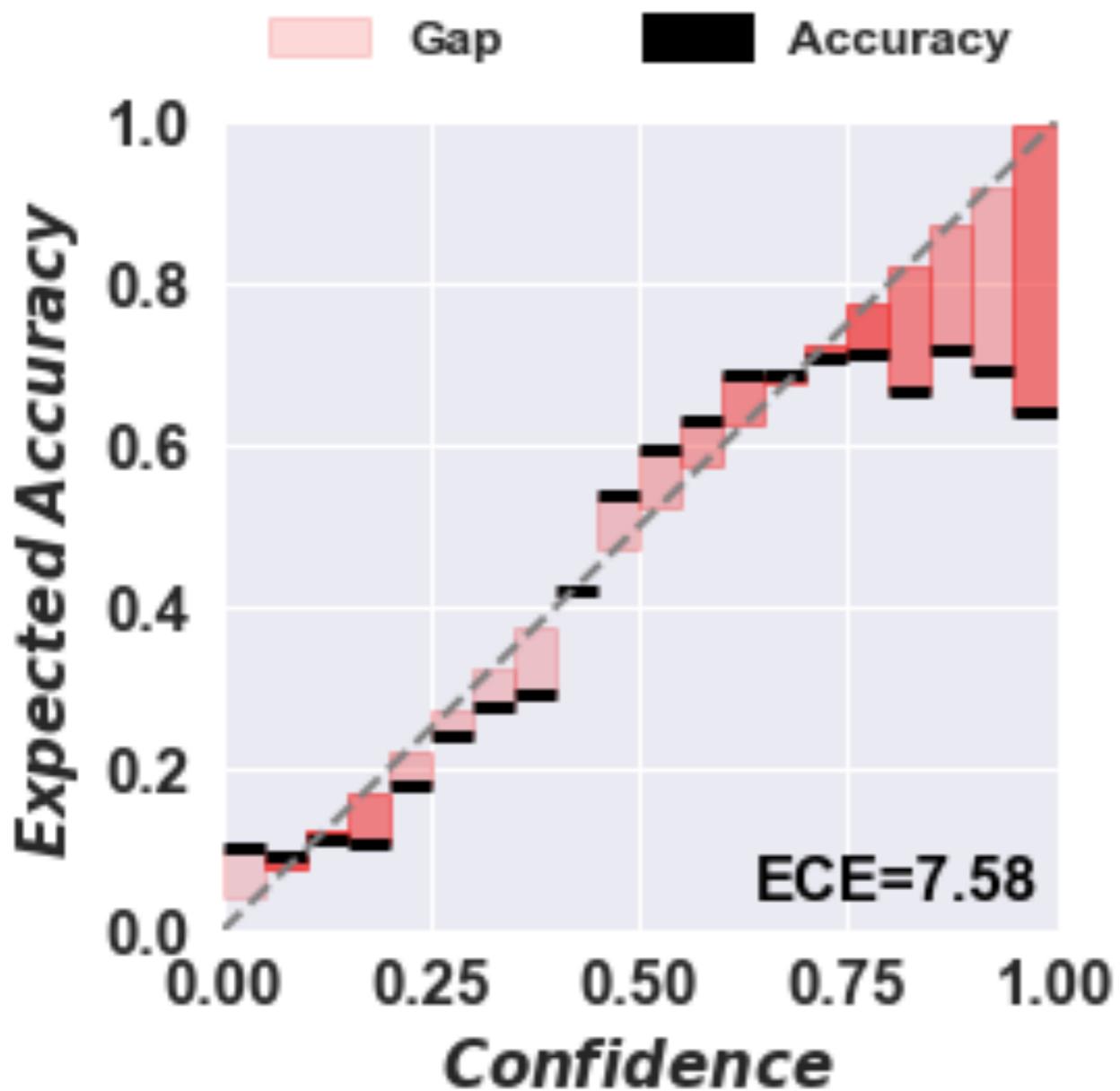
Theorem 4.1 + Corollary 4.2 [Verma & Nalisnick, 2022]:  
The one-vs-all loss is a *consistent* surrogate for  
the 0-1 learning-to-defer loss.

Proof: Via method of *error correcting output codes* [Dietterich & Bakiri, 1995], by applying  
Ramaswamy et al.'s [2018] Equation 1.

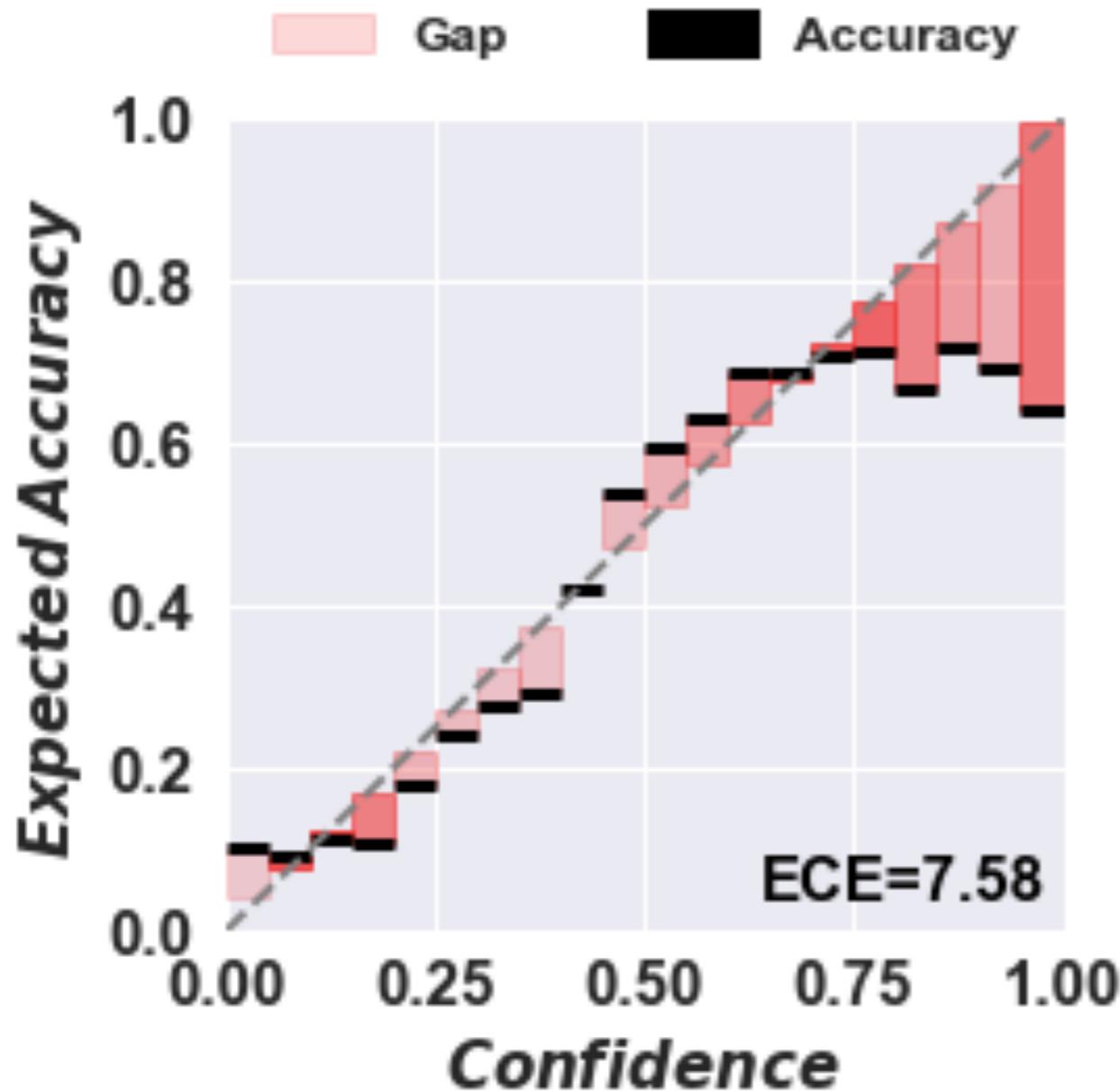
This means that our one-vs-all loss is as  
theoretically well-justified as Mozannar & Sontag's  
softmax-based surrogate.

Does the one-vs-all loss  
result in better calibrated  
models in practice?

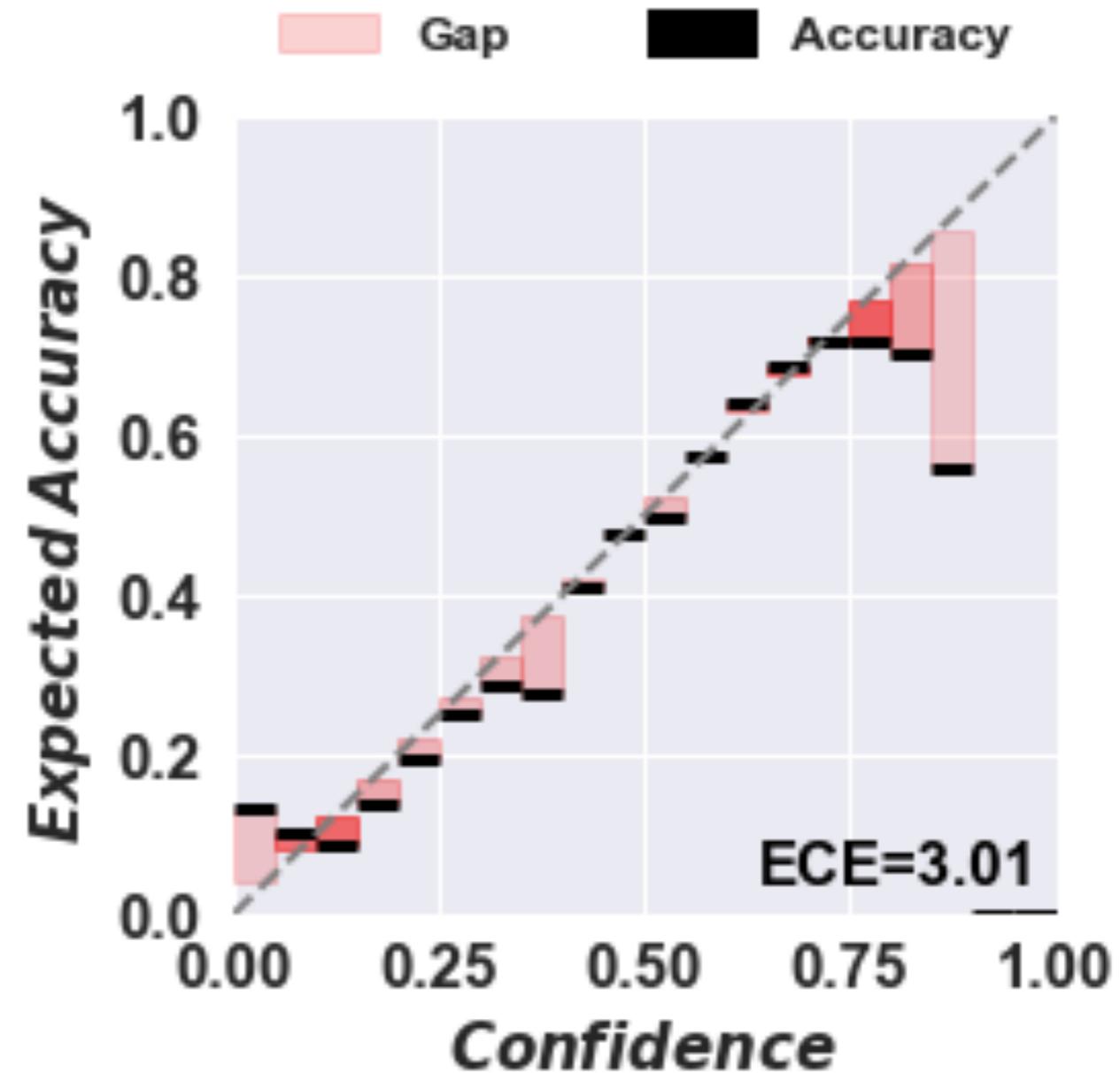
# CIFAR-10: Softmax



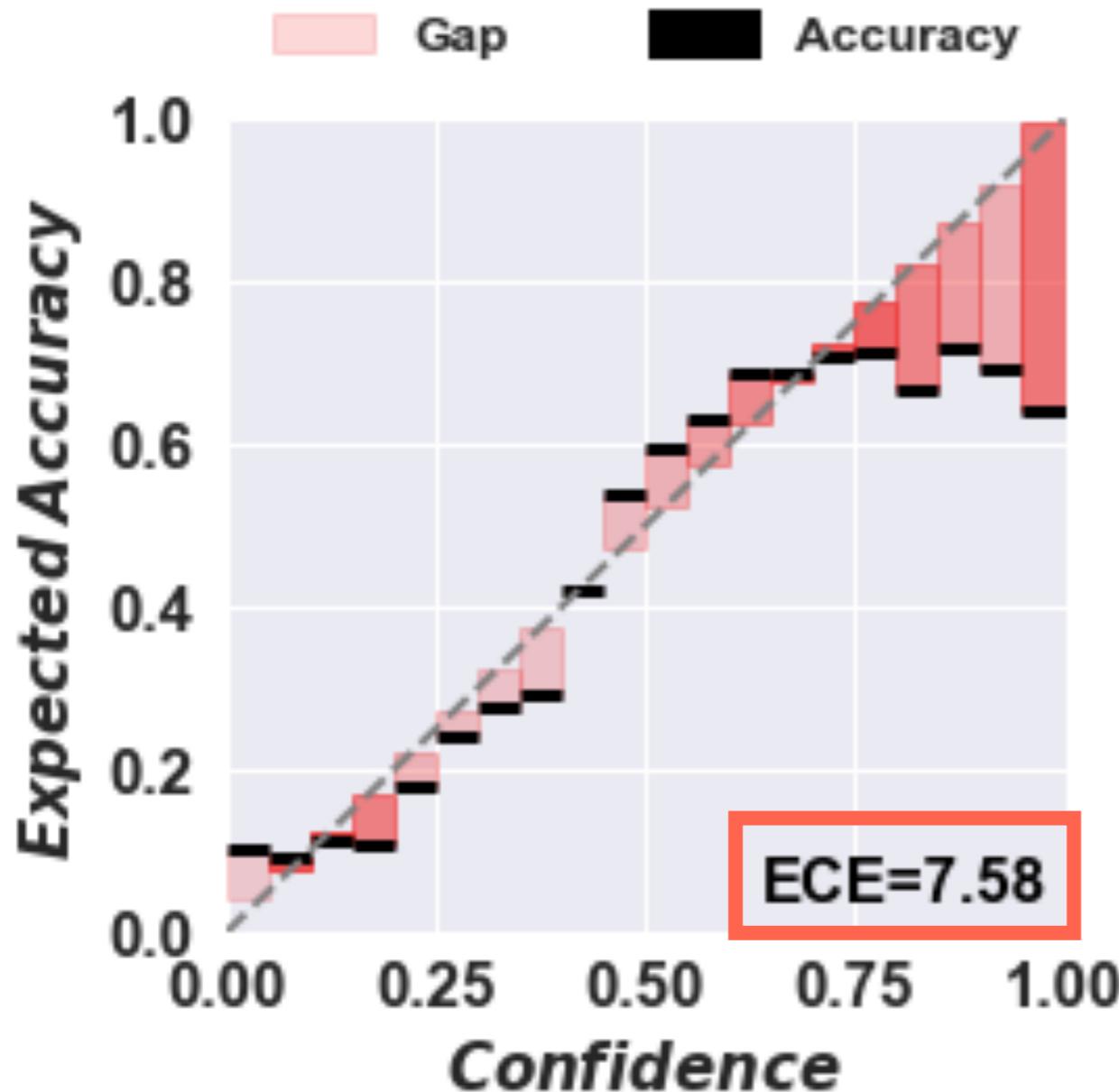
CIFAR-10: Softmax



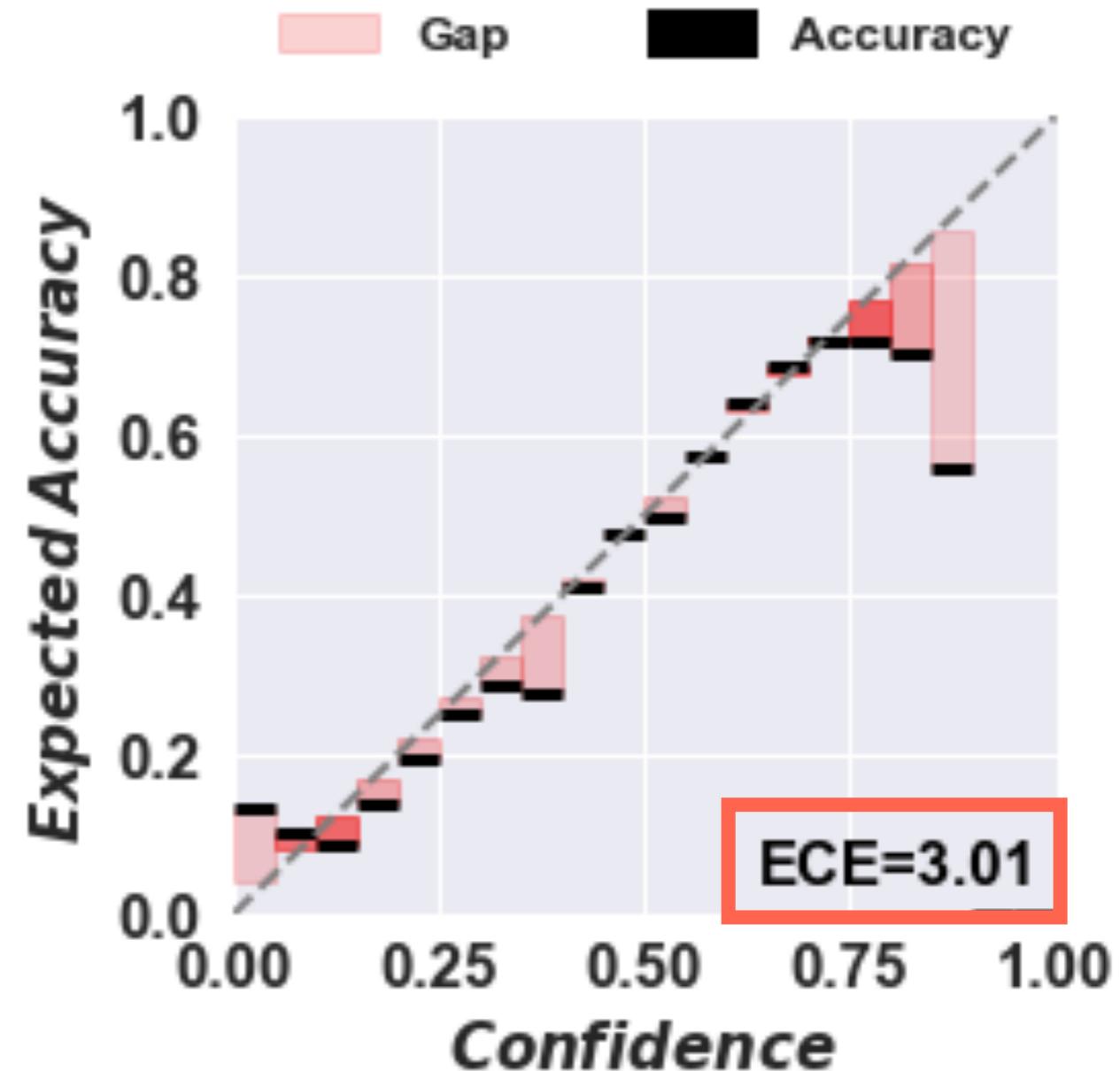
CIFAR-10: One-vs-All



CIFAR-10: Softmax



CIFAR-10: One-vs-All



# CIFAR-10: ECE Across Parameterizations

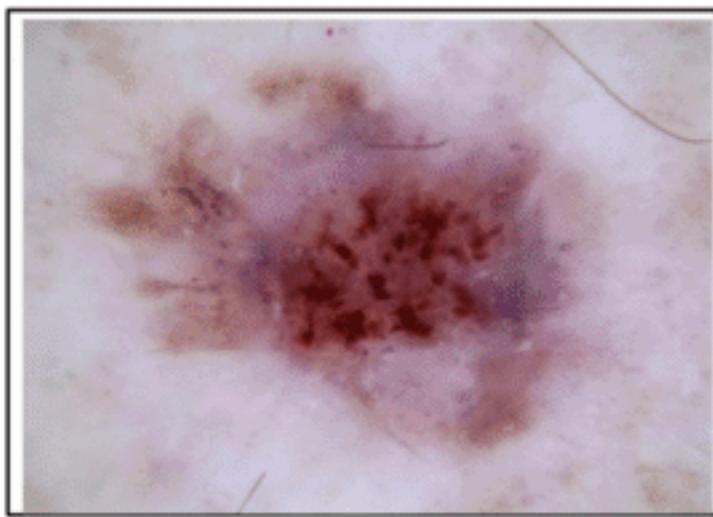
	OvA	Softmax	Proxy
Both Random			
Random Expert			
Random Data	$\sigma(g_{\perp}(x))$	$\frac{p_{\perp}(x)}{1 - p_{\perp}(x)}$	$p_{\perp}(x)$
Both Useful			

# CIFAR-10: ECE Across Parameterizations

	OvA	Softmax	Proxy
Both Random	0.53	0.97	<b>0.04</b>
Random Expert	<b>0.68</b>	3.72	2.83
Random Data	<b>2.05</b>	2.07	39.06
Both Useful	<b>1.68</b>	3.32	37.15

# HAM10000: Skin Lesion Classification

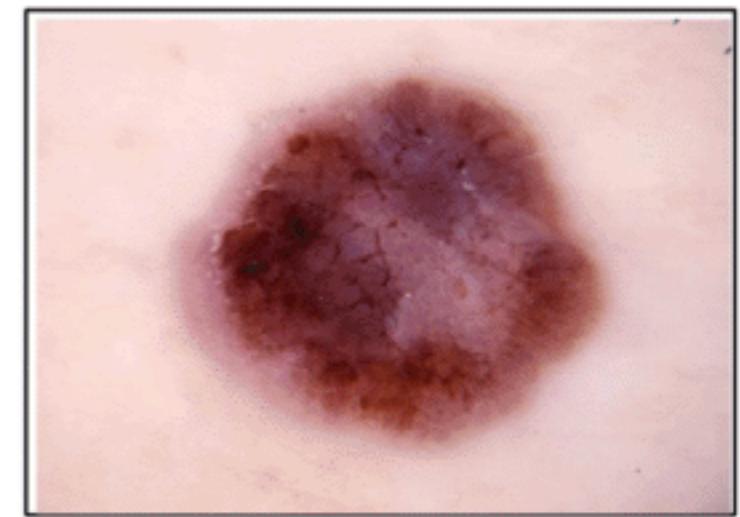
(AKIEC)



(BCC)



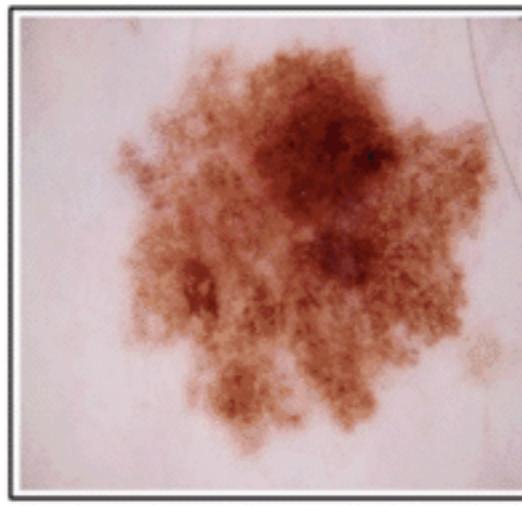
(BKL)



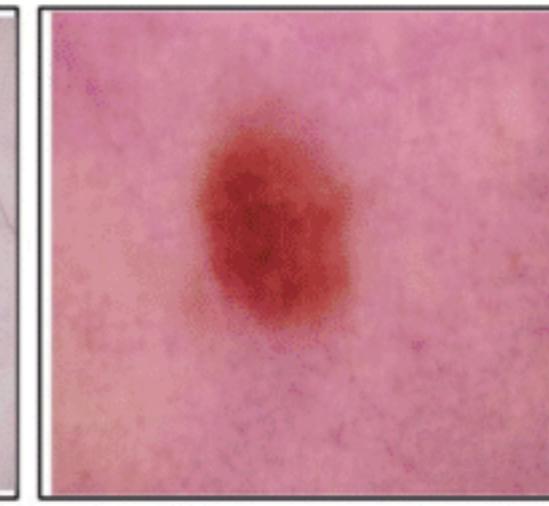
(DF)



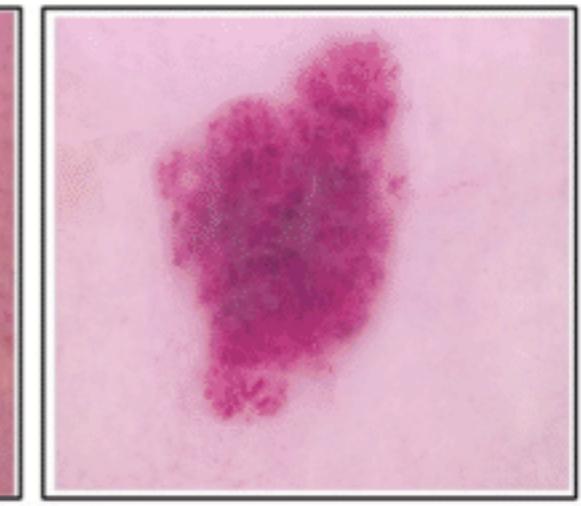
(MEL)



(NV)



(VASC)



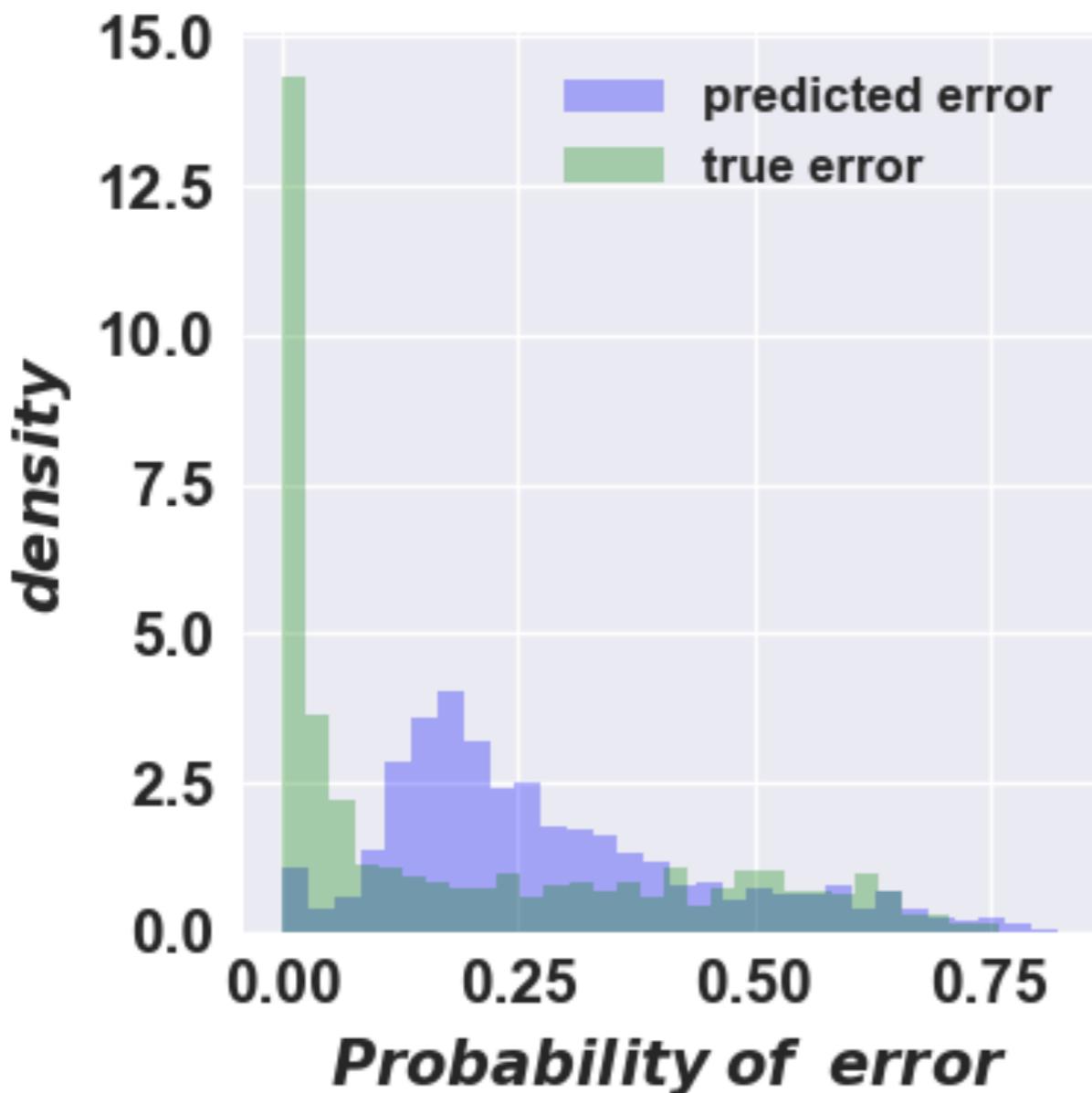
# HAM10000: Skin Lesion Classification

Softmax

One-vs-All

# HAM10000: Skin Lesion Classification

Softmax

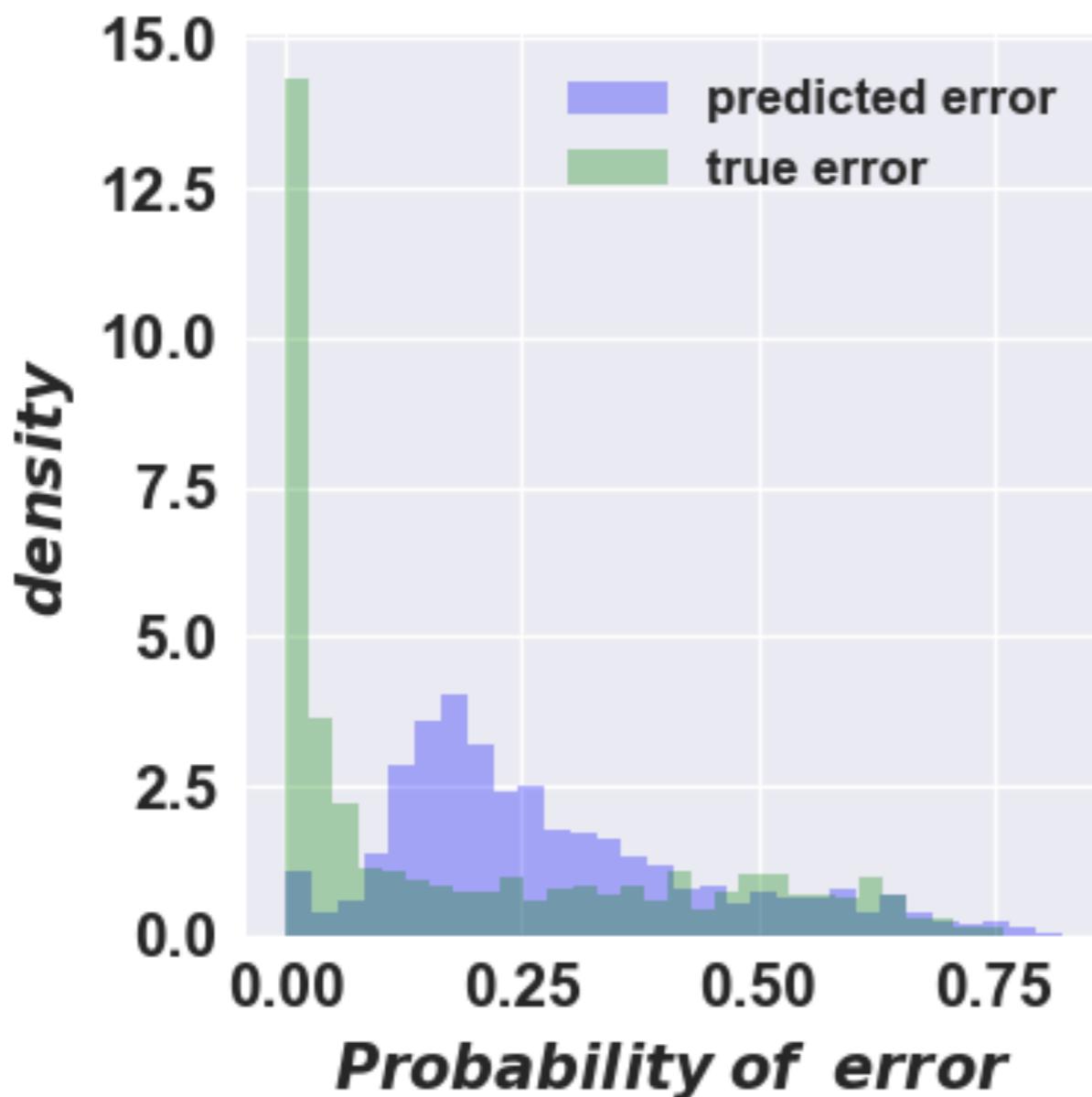


One-vs-All

$$1 - p_m(x)$$

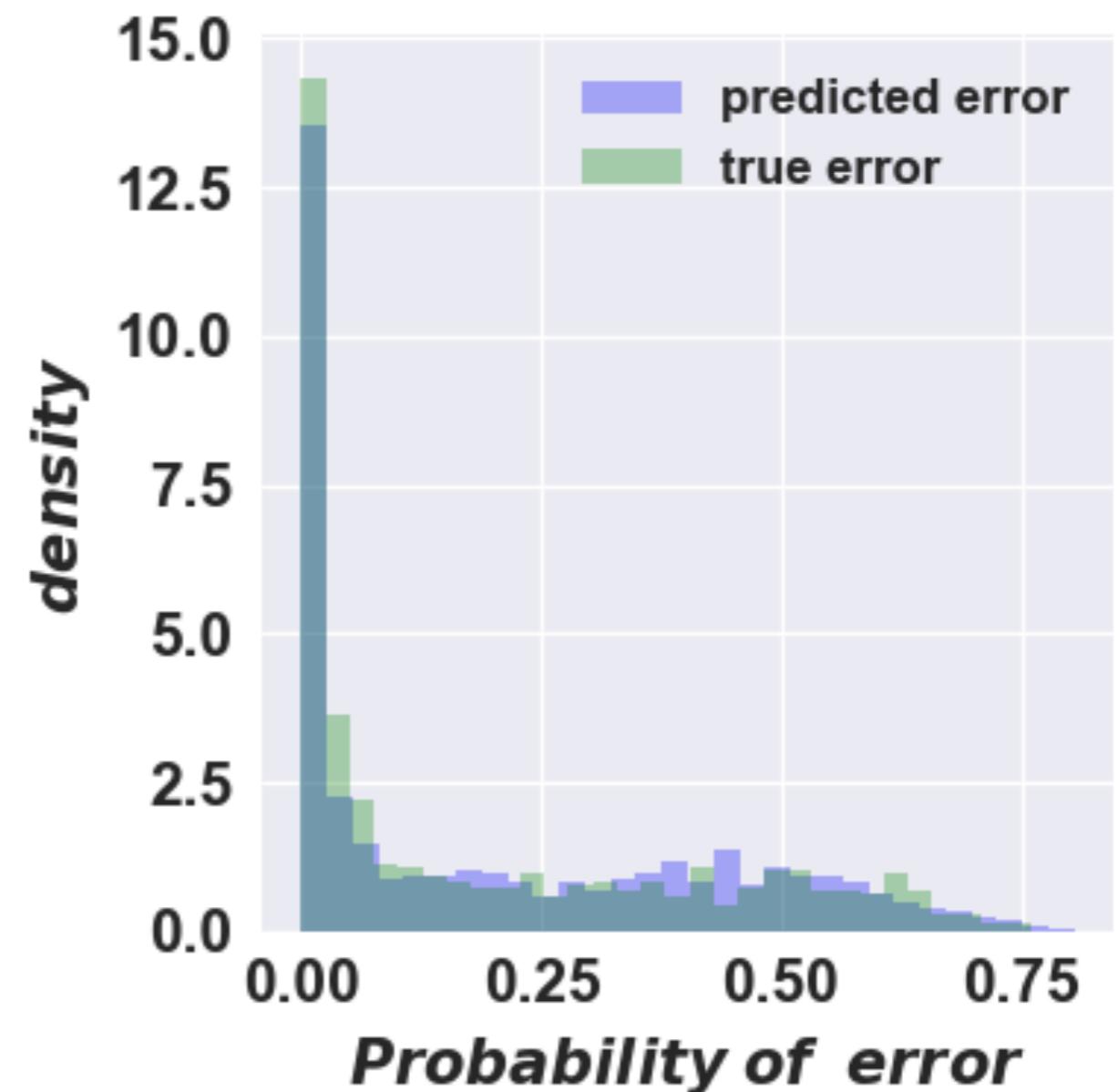
# HAM10000: Skin Lesion Classification

Softmax



$$1 - p_m(x)$$

One-vs-All

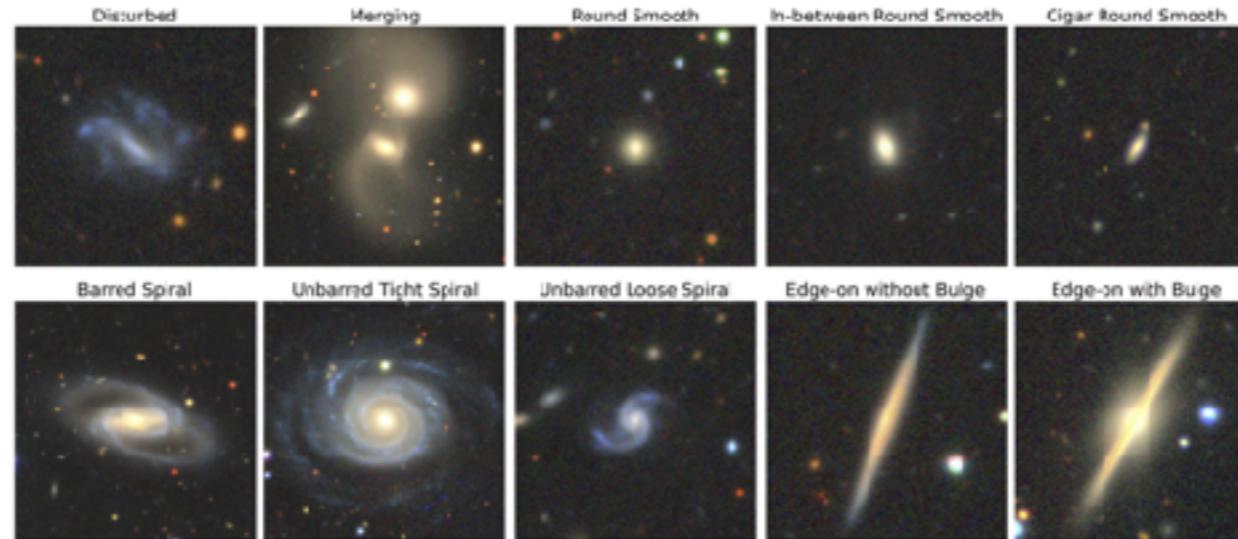


$$1 - p_m^{\text{OvA}}(x)$$

Does the one-vs-all loss  
result in more accurate  
models in practice?

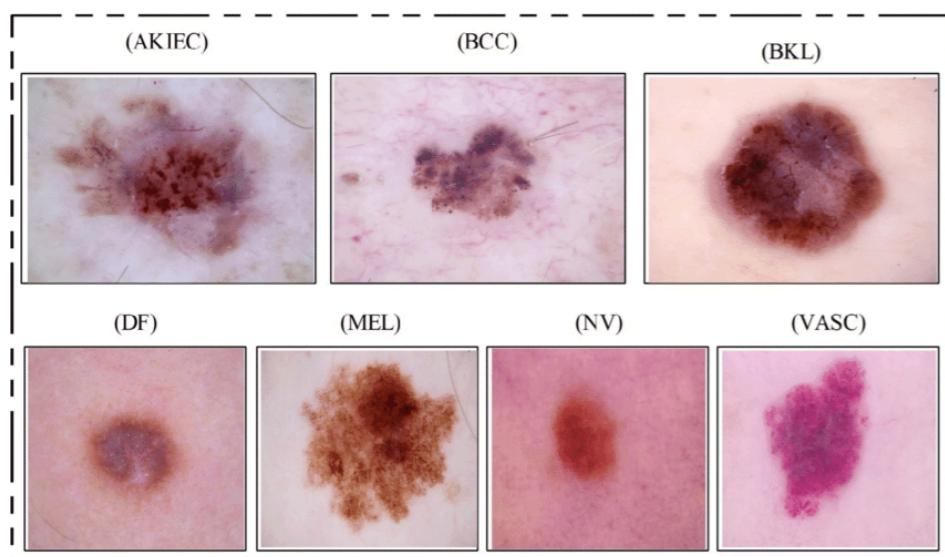
one-vs-all (ours) softmax confidence score differentiable triage

# Galaxy Zoo



Galaxy10 DECaLS: Henry Leung/o Bovy 2021, Data: DECaLS/Galaxy Zoo

# HAM10000



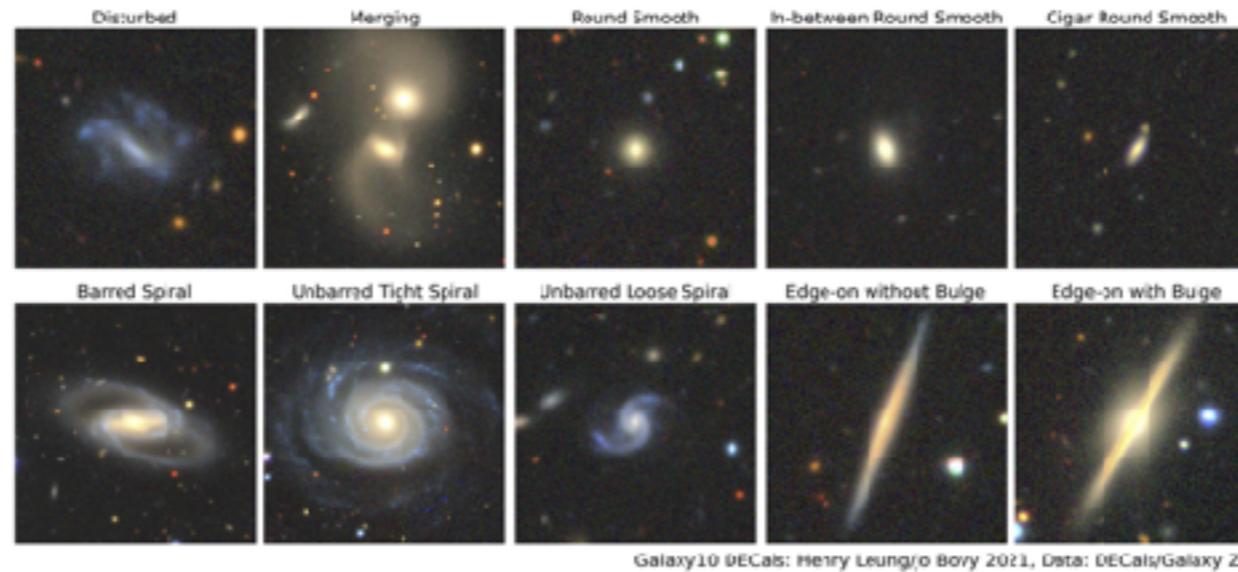
# Hate Speech



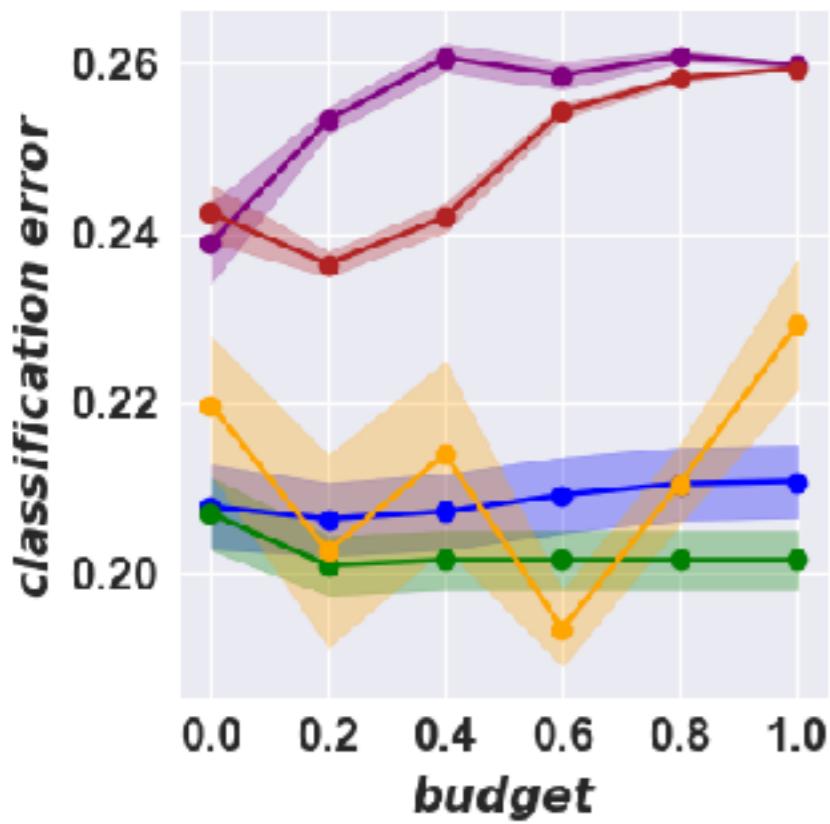
[Davidson et al., ICWSM 2017]

one-vs-all (ours) softmax confidence score differentiable triage

## Galaxy Zoo



## HAM10000



## Hate Speech



[Davidson et al., ICWSM 2017]

one-vs-all (ours)

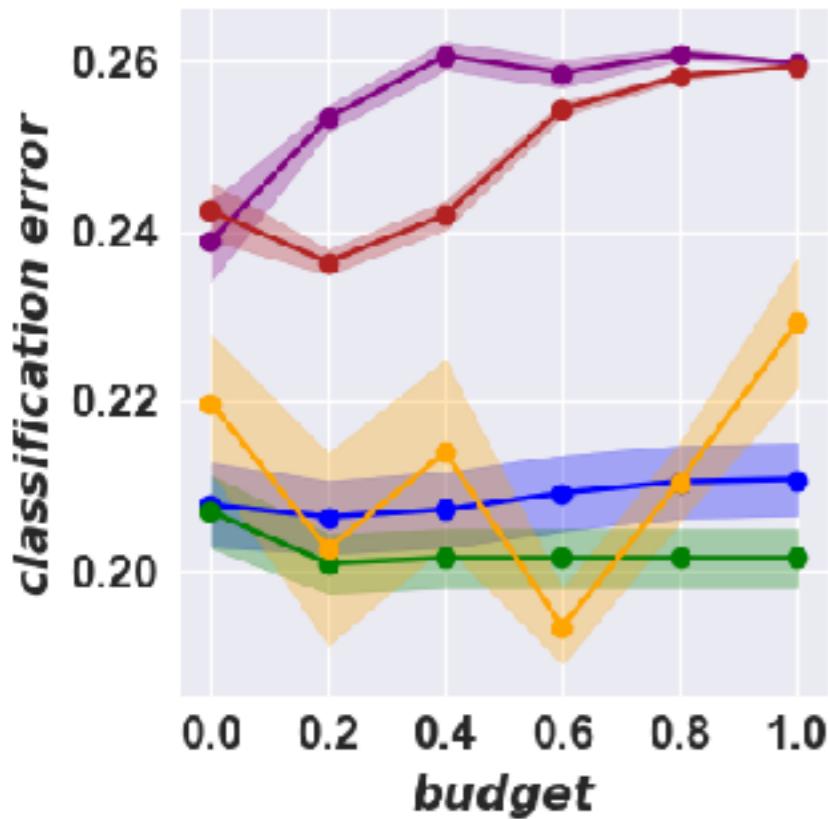
softmax

confidence

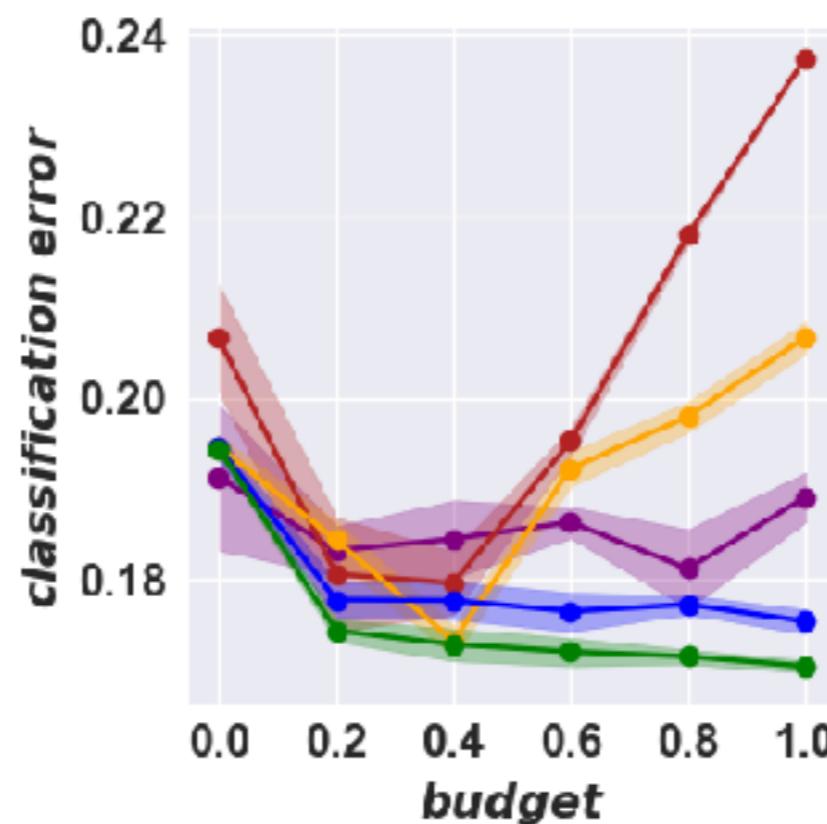
score

differentiable triage

# HAM10000



# Galaxy Zoo



# Hate Speech



[Davidson et al., ICWSM 2017]

one-vs-all (ours)

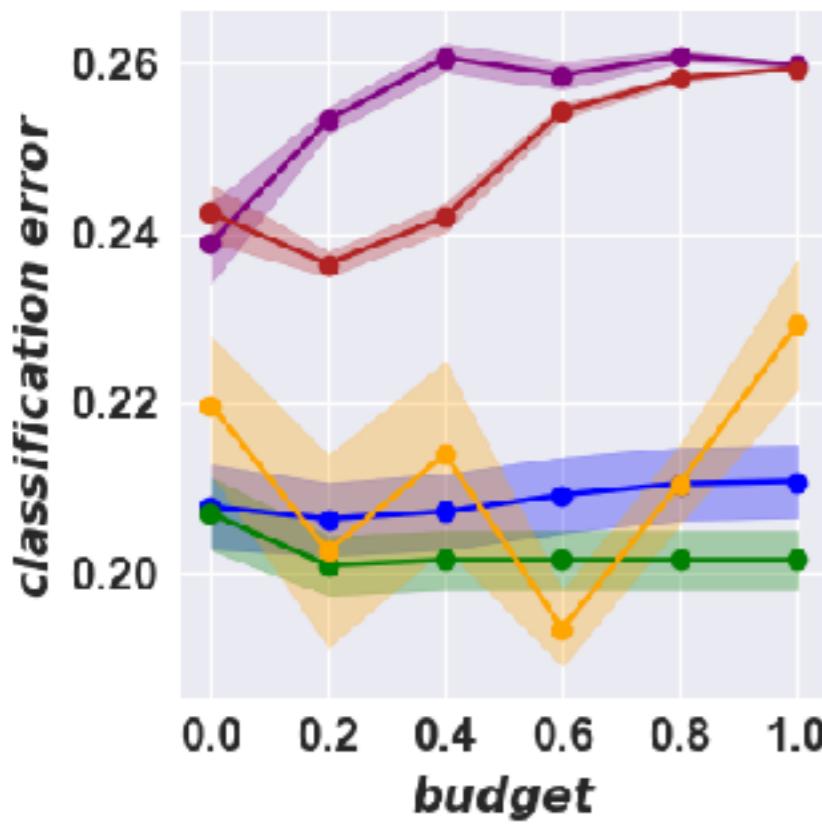
softmax

confidence

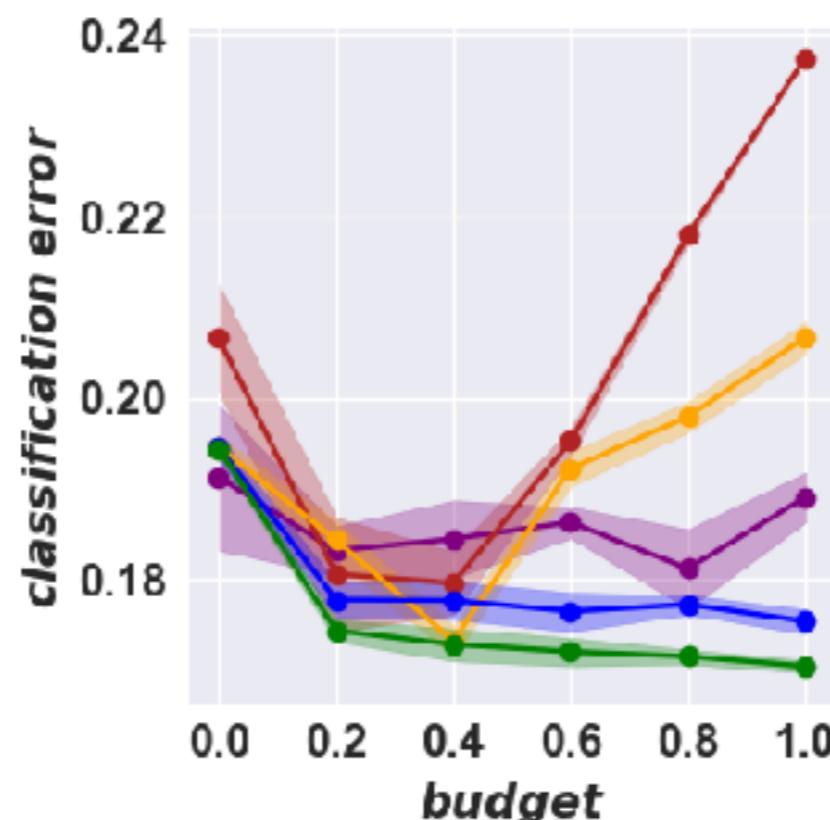
score

differentiable triage

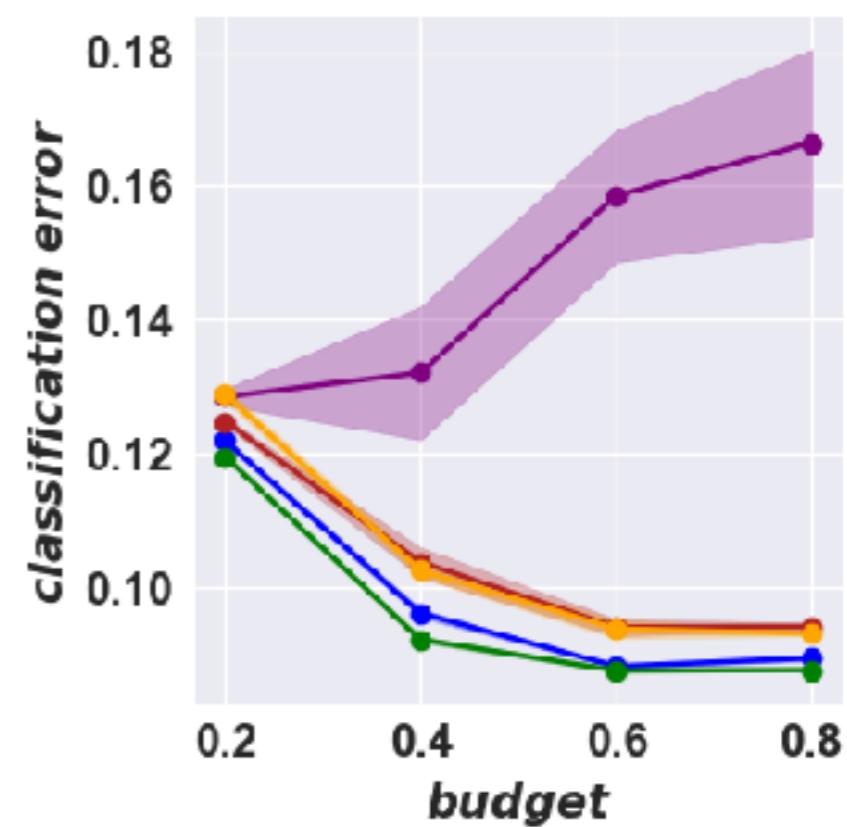
# HAM10000



# Galaxy Zoo

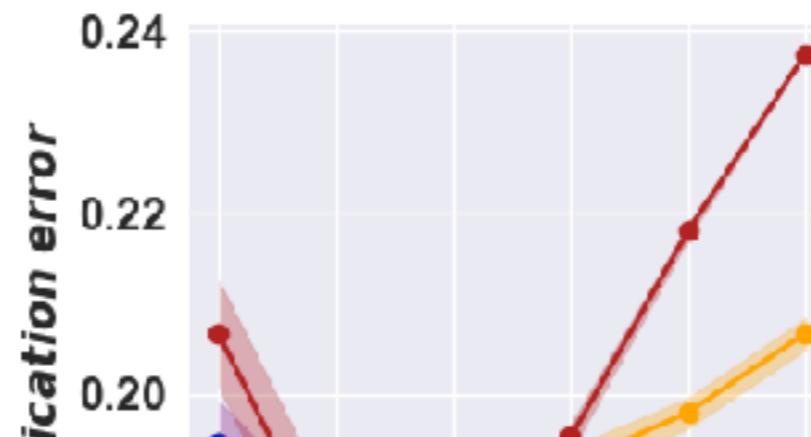


# Hate Speech

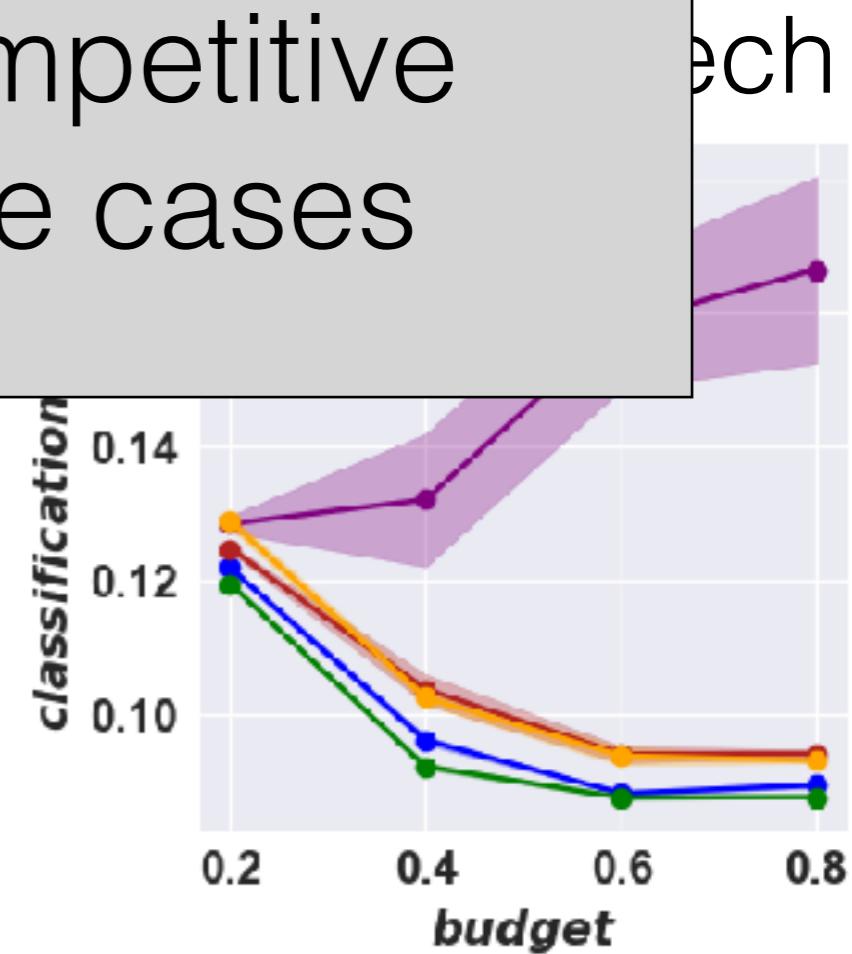
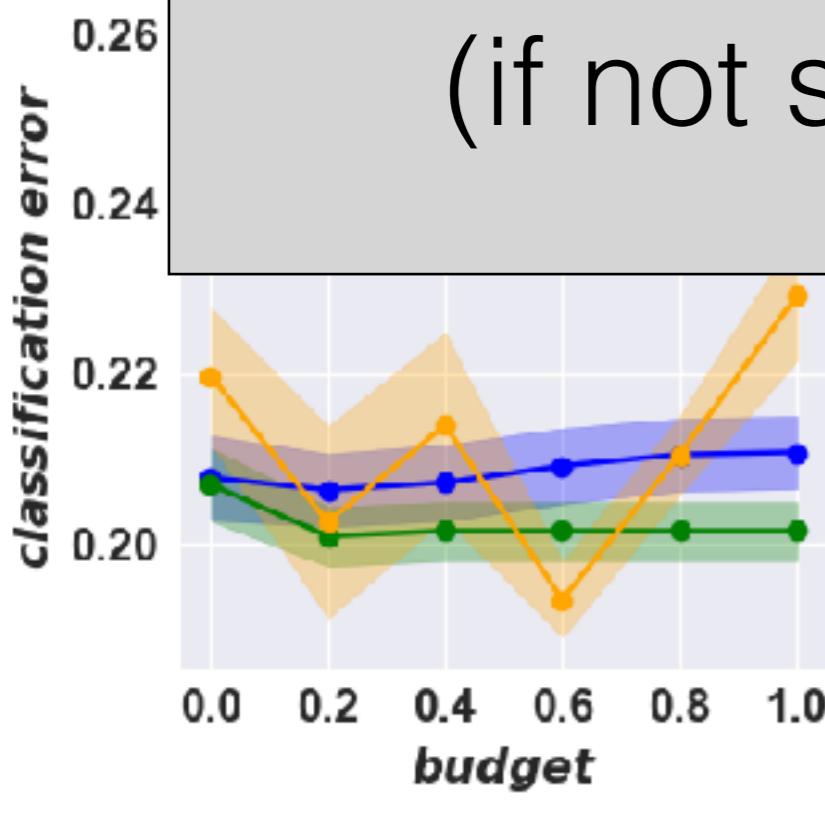


one-vs-all (ours) softmax confidence score differentiable triage

# Galaxy Zoo



One-vs-all models are competitive  
(if not superior) in all three cases



# Summary

- ⊗ We show that the softmax-based loss for learning-to-defer can produce invalid estimators for expert correctness.
- ⊗ Our one-vs-all formulation retains theoretical and practical benefits while also producing calibrated estimates of expert correctness.

# Future Work

- ⊗ Learning-to-defer is still under-studied!
- ⊗ One problem: heavy dependence on supervised data.

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


## Computer Science &gt; Machine Learning

*[Submitted on 8 Feb 2022]*

# Calibrated Learning to Defer with One-vs-All Classifiers

[Rajeev Verma, Eric Nalisnick](#)

The learning to defer (L2D) framework has the potential to make AI systems safer. For a given input, the system can defer the decision to a human if the human is more likely than the model to take the correct action. We study the calibration of L2D systems, investigating if the probabilities they output are sound. We find that Mozannar & Sontag's (2020) multiclass framework is not calibrated with respect to expert correctness. Moreover, it is not even guaranteed to produce valid probabilities due to its parameterization being degenerate for this purpose. We propose an L2D system based on one-vs-all classifiers that is able to produce calibrated probabilities of expert correctness. Furthermore, our loss function is also a consistent surrogate for multiclass L2D, like Mozannar & Sontag's (2020). Our experiments verify that not only is our system calibrated, but this benefit comes at no cost to accuracy. Our model's accuracy is always comparable (and often superior) to Mozannar & Sontag's (2020) model's in tasks ranging from hate speech detection to galaxy classification to diagnosis of skin lesions.

Subjects: [Machine Learning \(cs.LG\)](#); [Machine Learning \(stat.ML\)](#)

Cite as: [arXiv:2202.03673 \[cs.LG\]](#)

(or [arXiv:2202.03673v1 \[cs.LG\]](#) for this version)

Thank you. Questions?