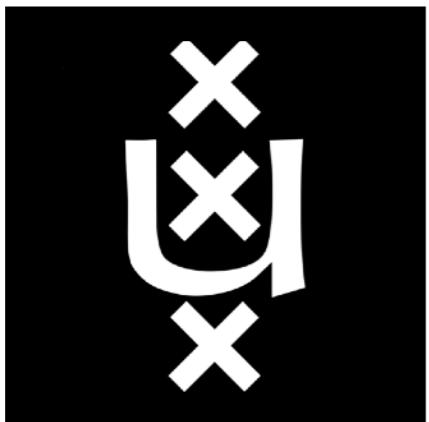

On the Calibration of Learning-to-Defer Systems

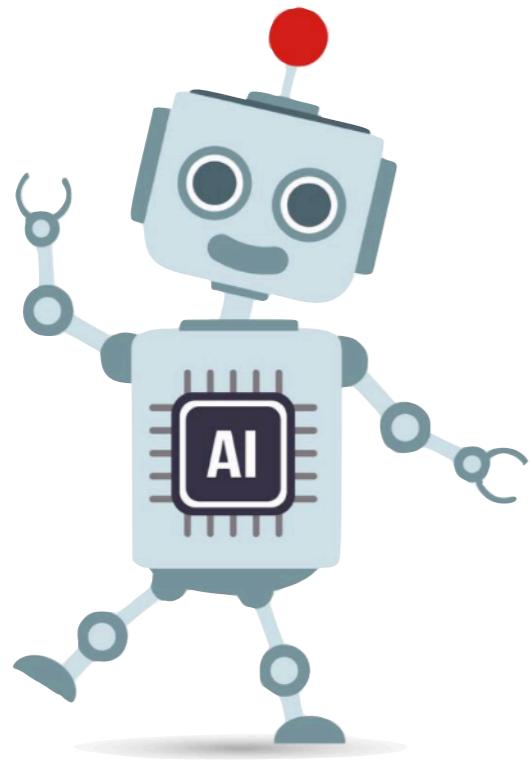
Eric Nalisnick



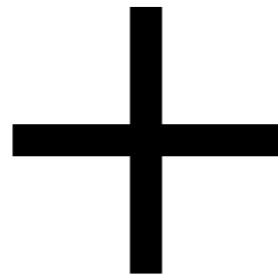
Rajeev
Verma



Daniel
Barrejón



Model

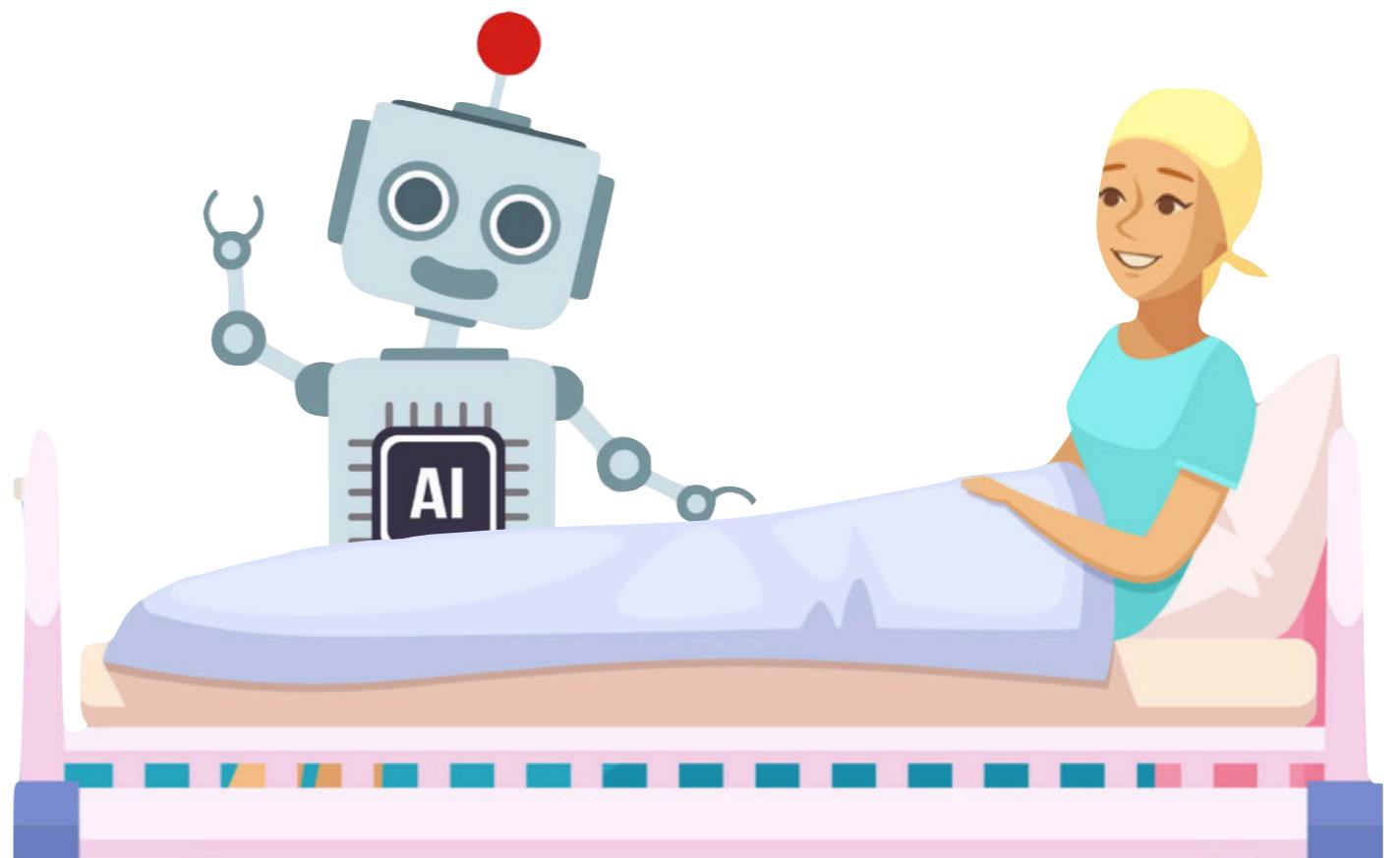


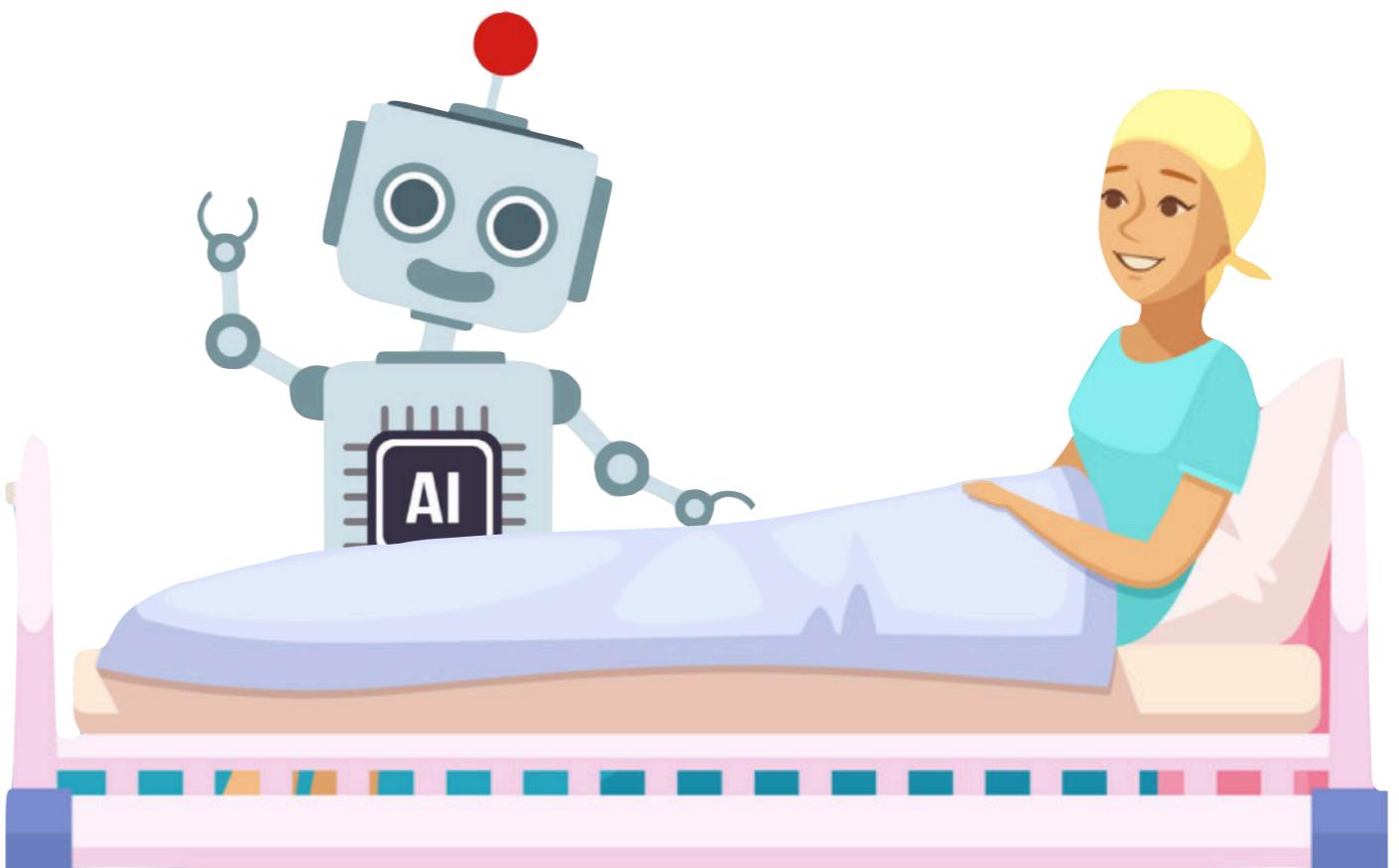
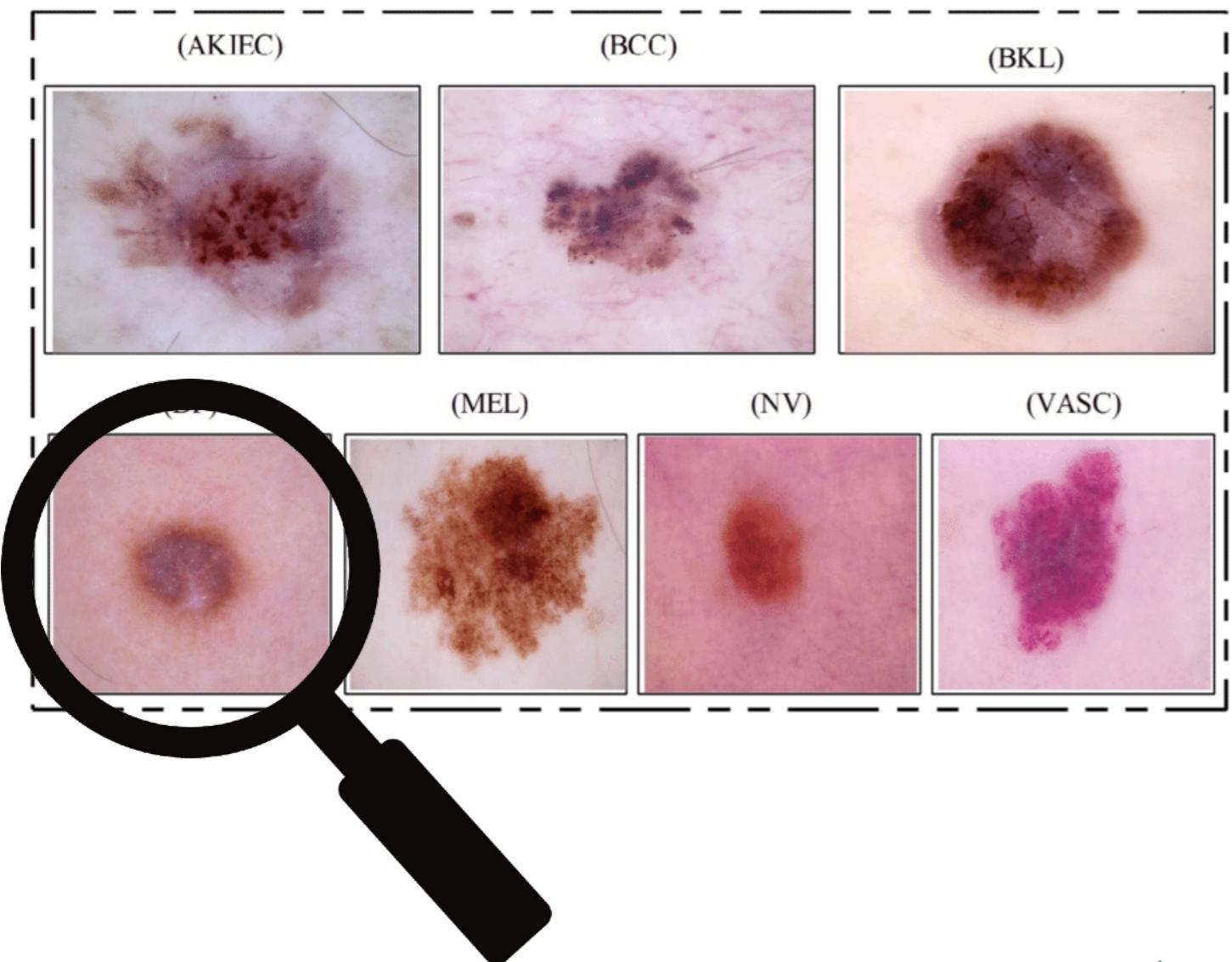
Human Expert

Artificial Intelligence

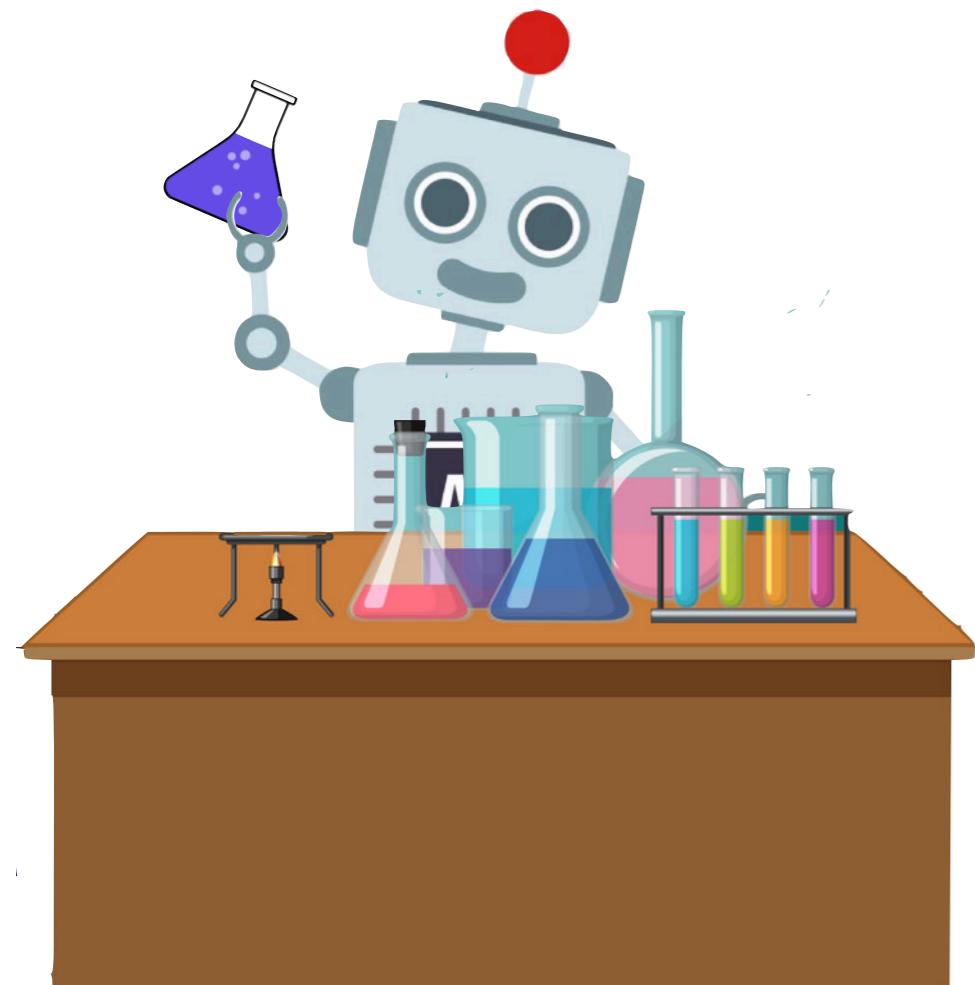
Hybrid Intelligence

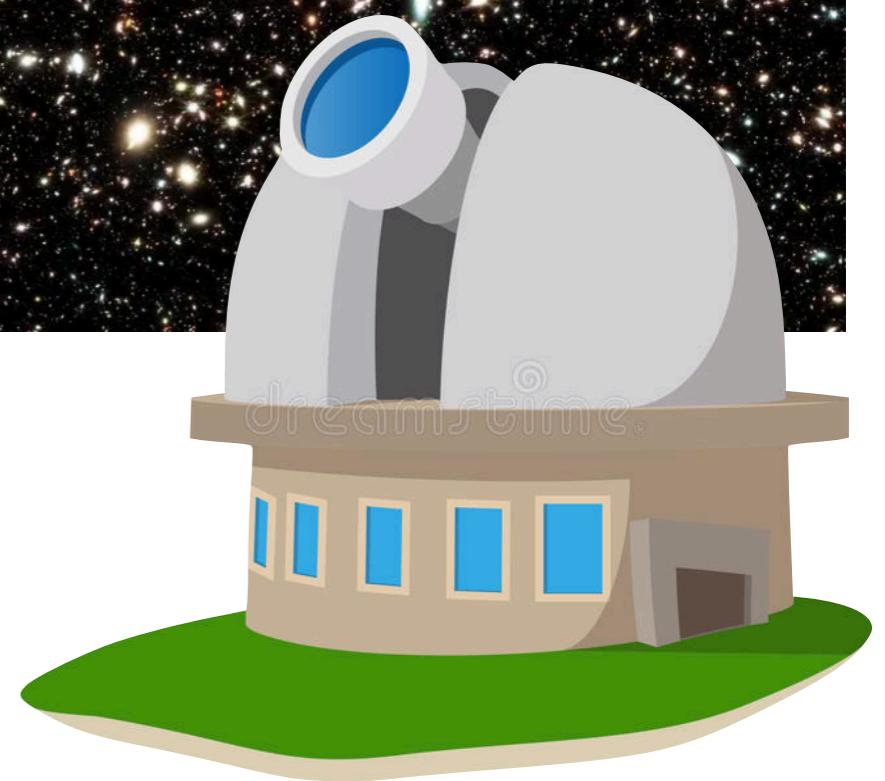
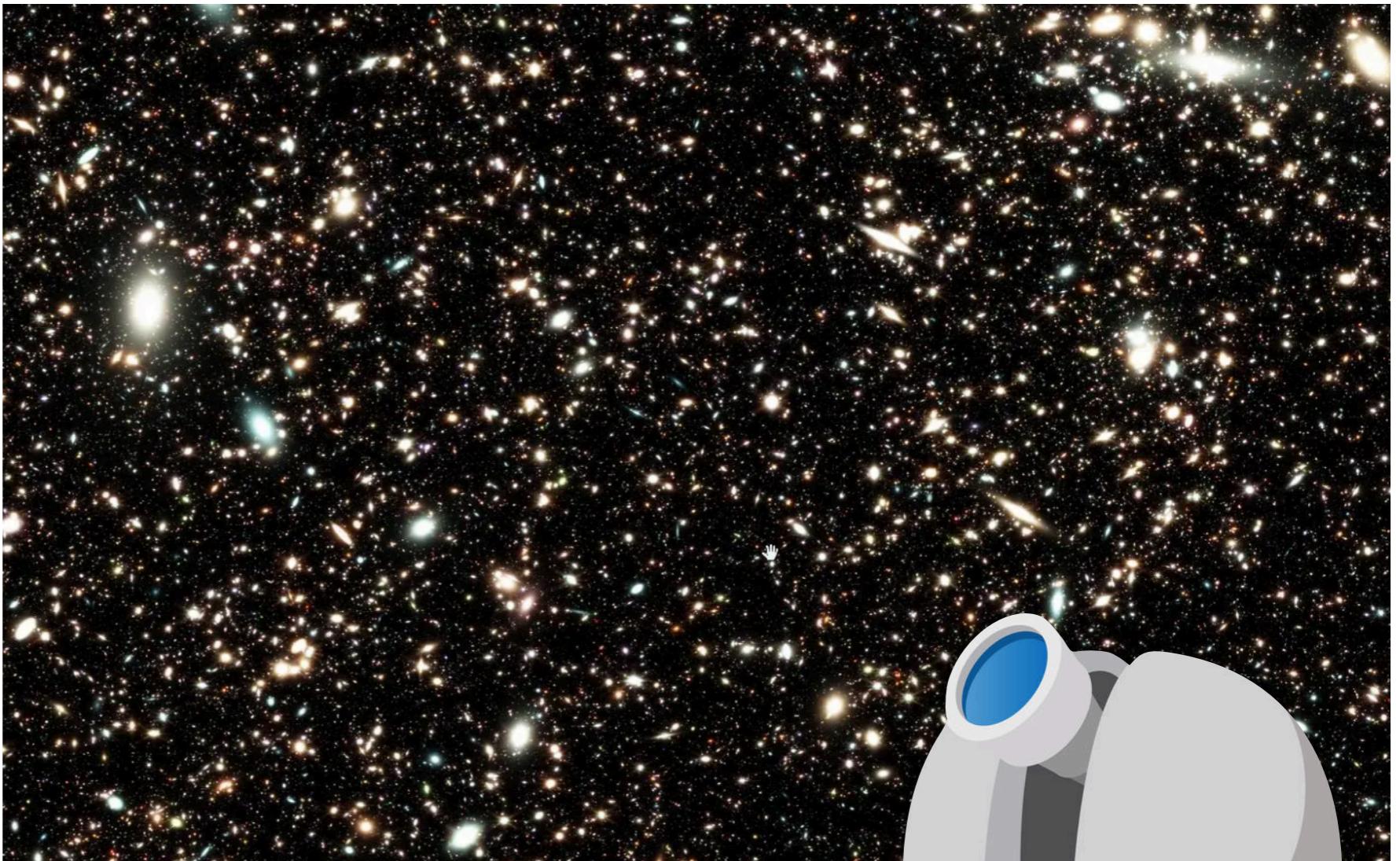
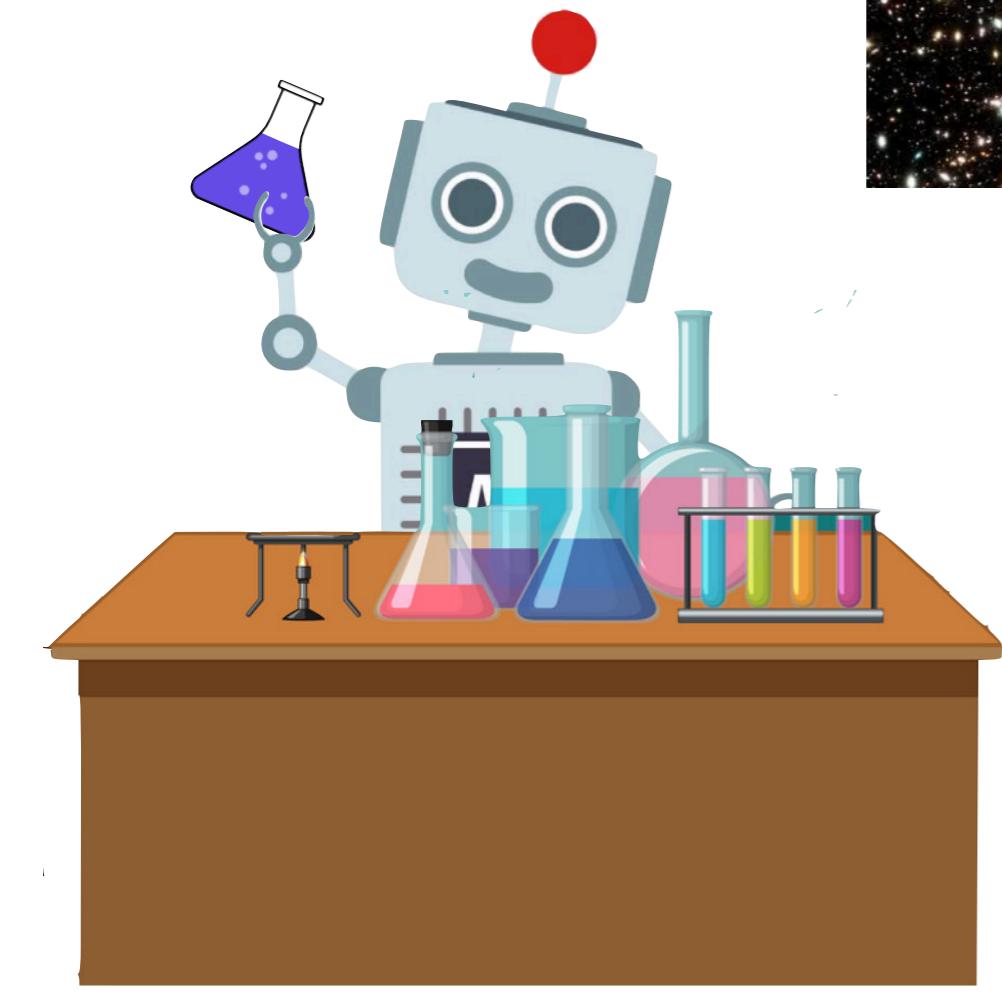
Safe Automation in Healthcare



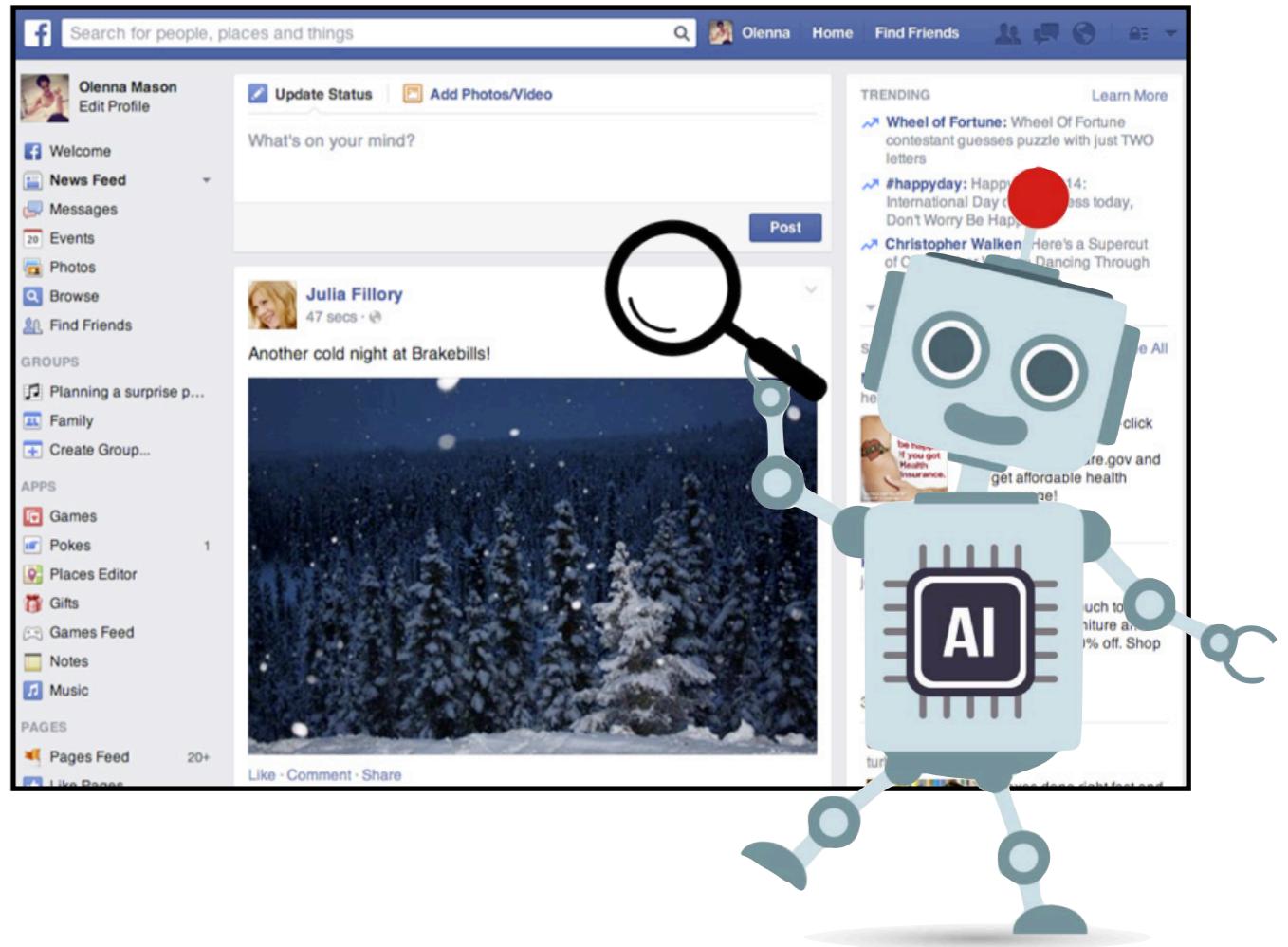


Large-Scale Science





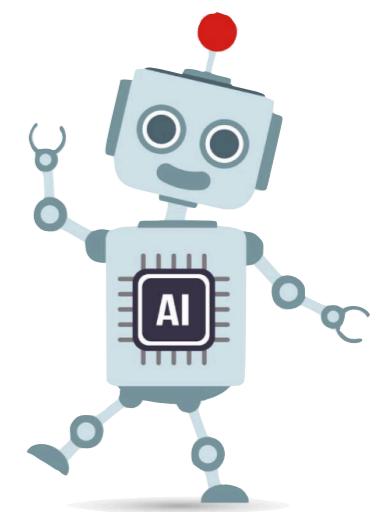
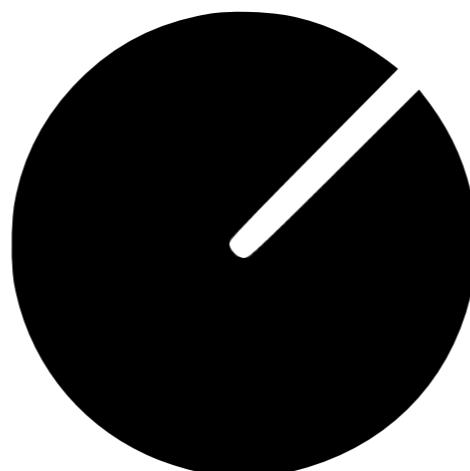
Online Content Moderation

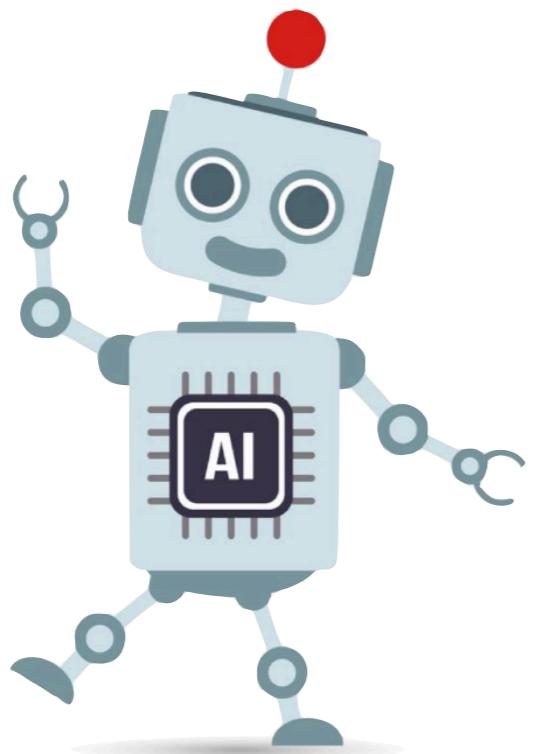




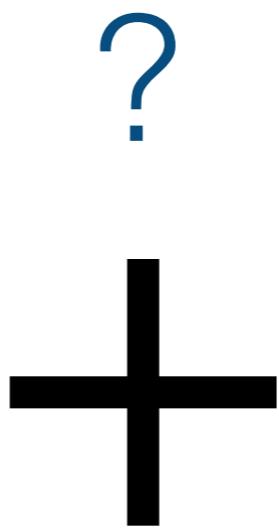
A screenshot of a Facebook news feed. On the left, there's a sidebar with links like Welcome, News Feed, Messages, Events, Photos, Browse, Find Friends, Groups, Planning a surprise p..., Family, Create Group..., Games, Pokes, Places Editor, Gifts, Games Feed, Notes, Music, Pages Feed, and Like Pages. The main area shows a post from 'Julia Fillory' with the caption 'Another cold night at Brakebills!' and a photo of snow-covered trees. A magnifying glass icon is positioned over the photo. To the right, there's a 'Trending' section with items like 'Wheel of Fortune', '#happyday', and 'Christopher Walken'. A large, stylized blue robot character with the letters 'AI' on its chest is standing on the right side of the screen, holding a magnifying glass over the Facebook feed.

Gradual Automation

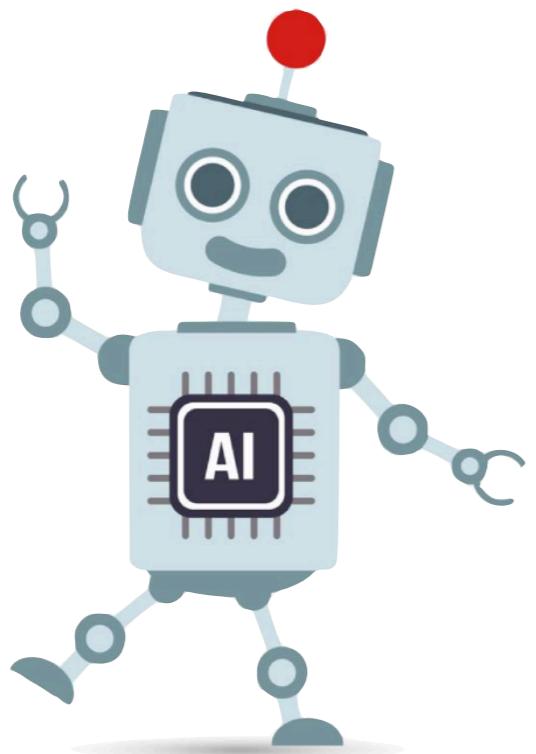




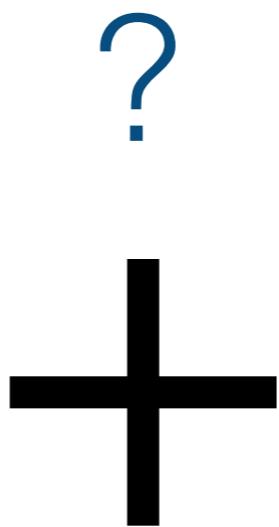
Model



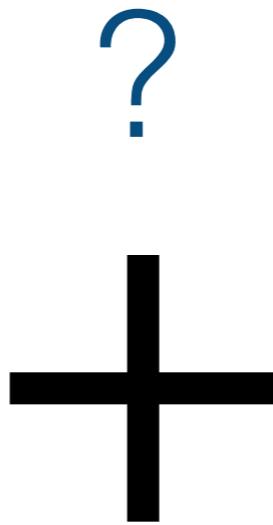
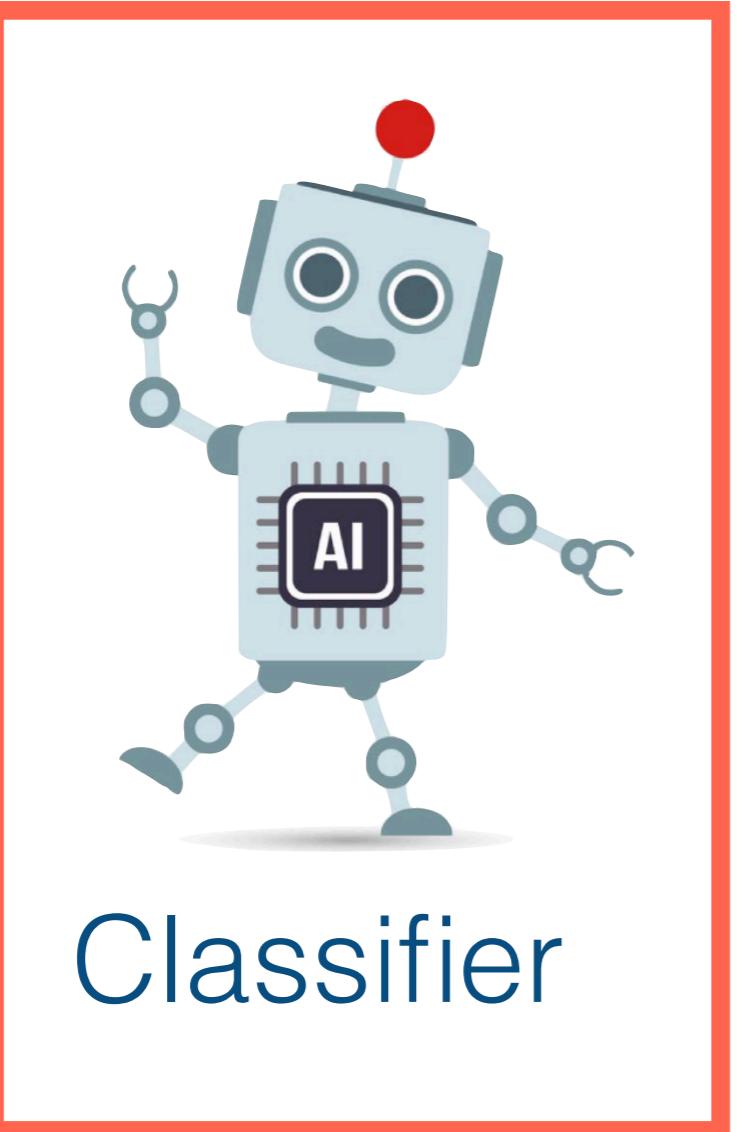
Human Expert



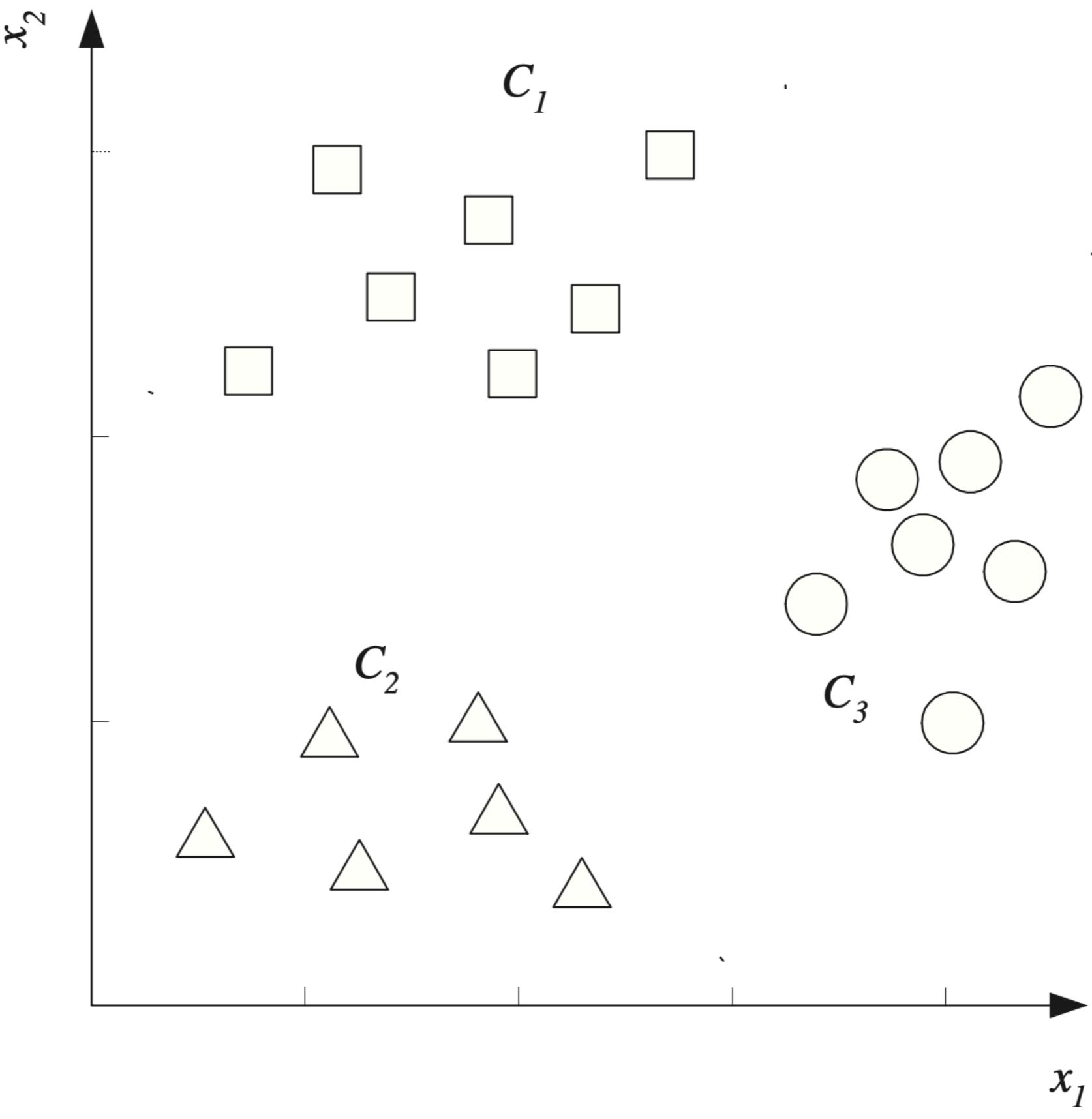
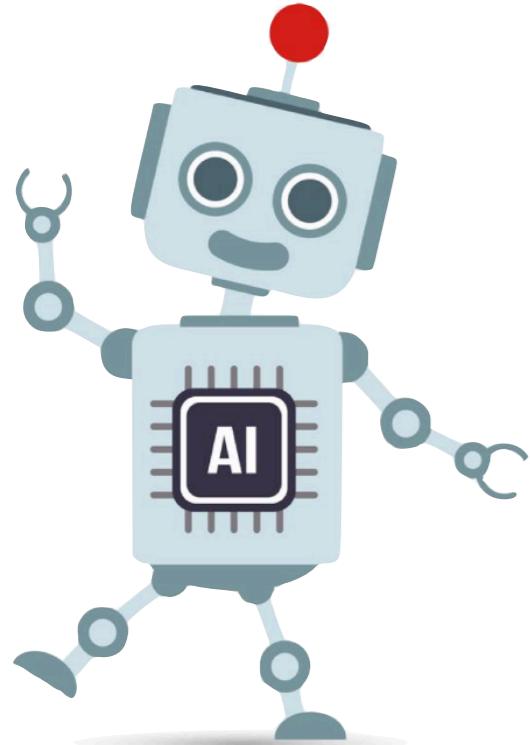
Classifier



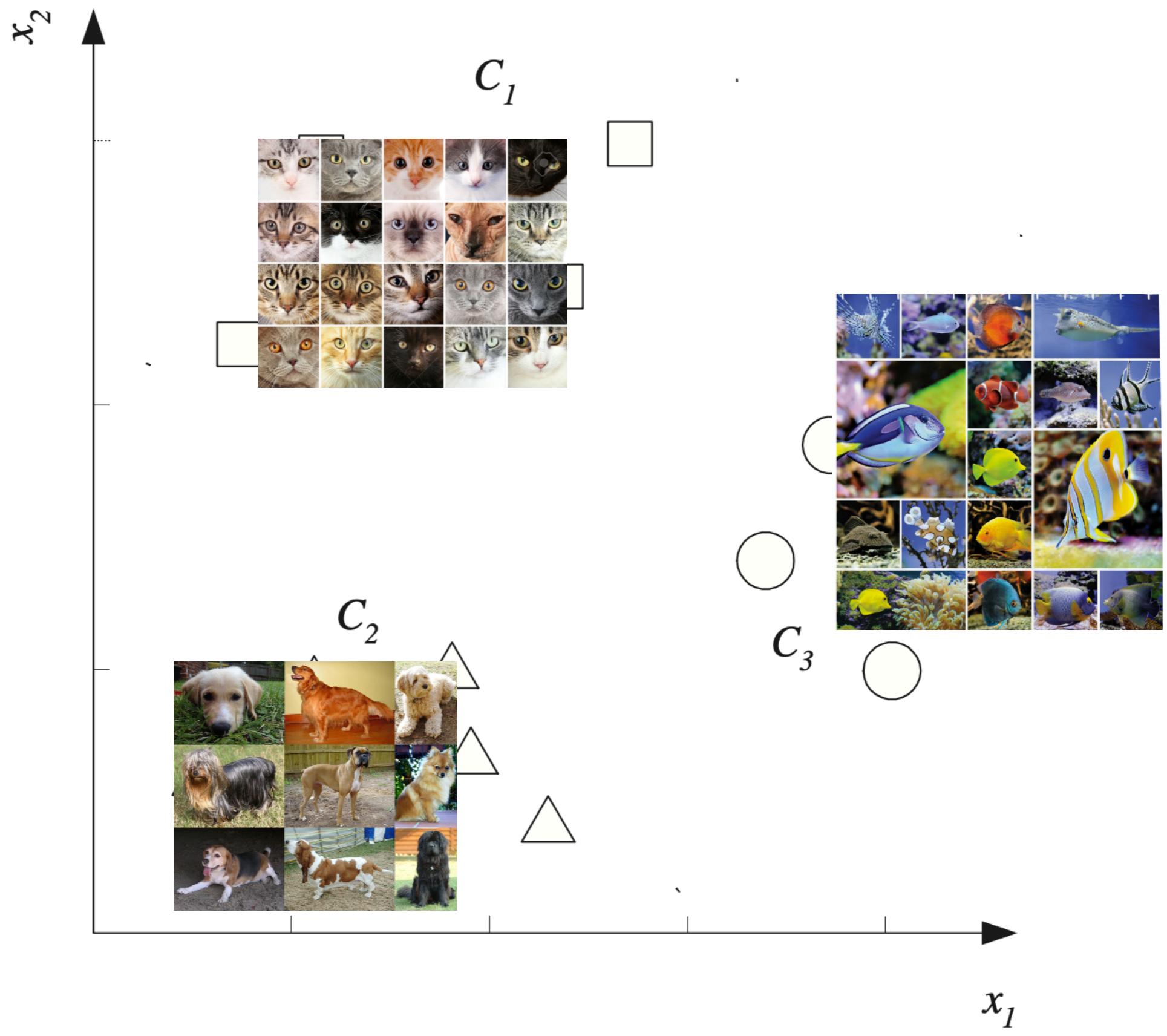
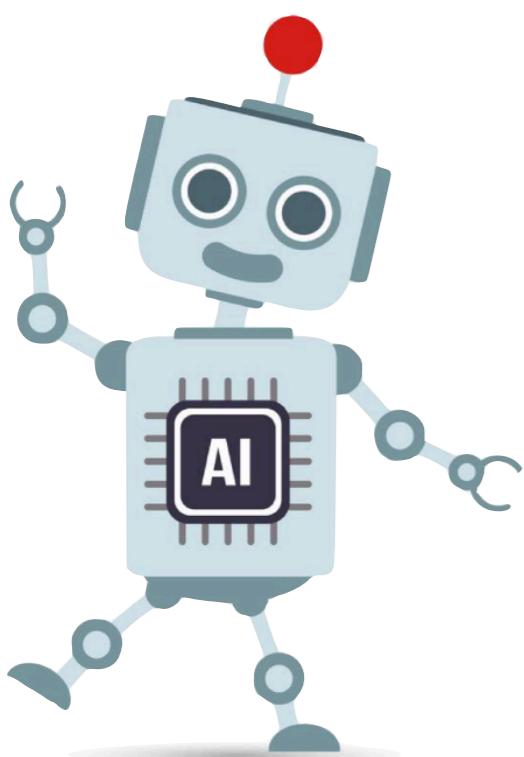
Human Expert



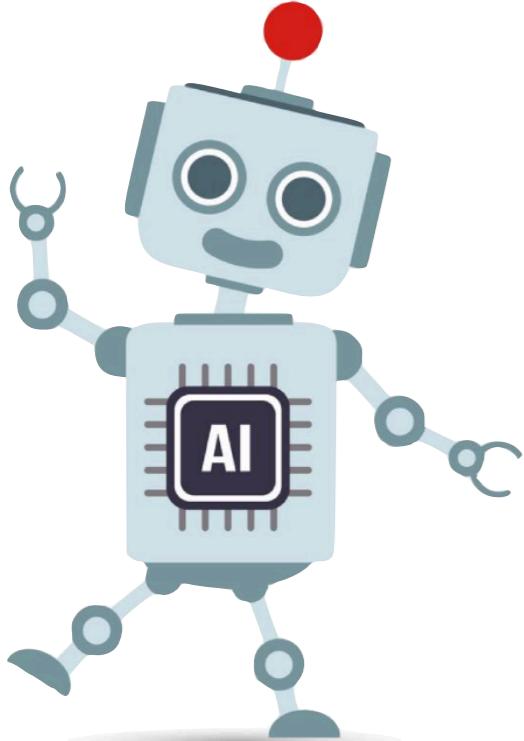
Human Expert



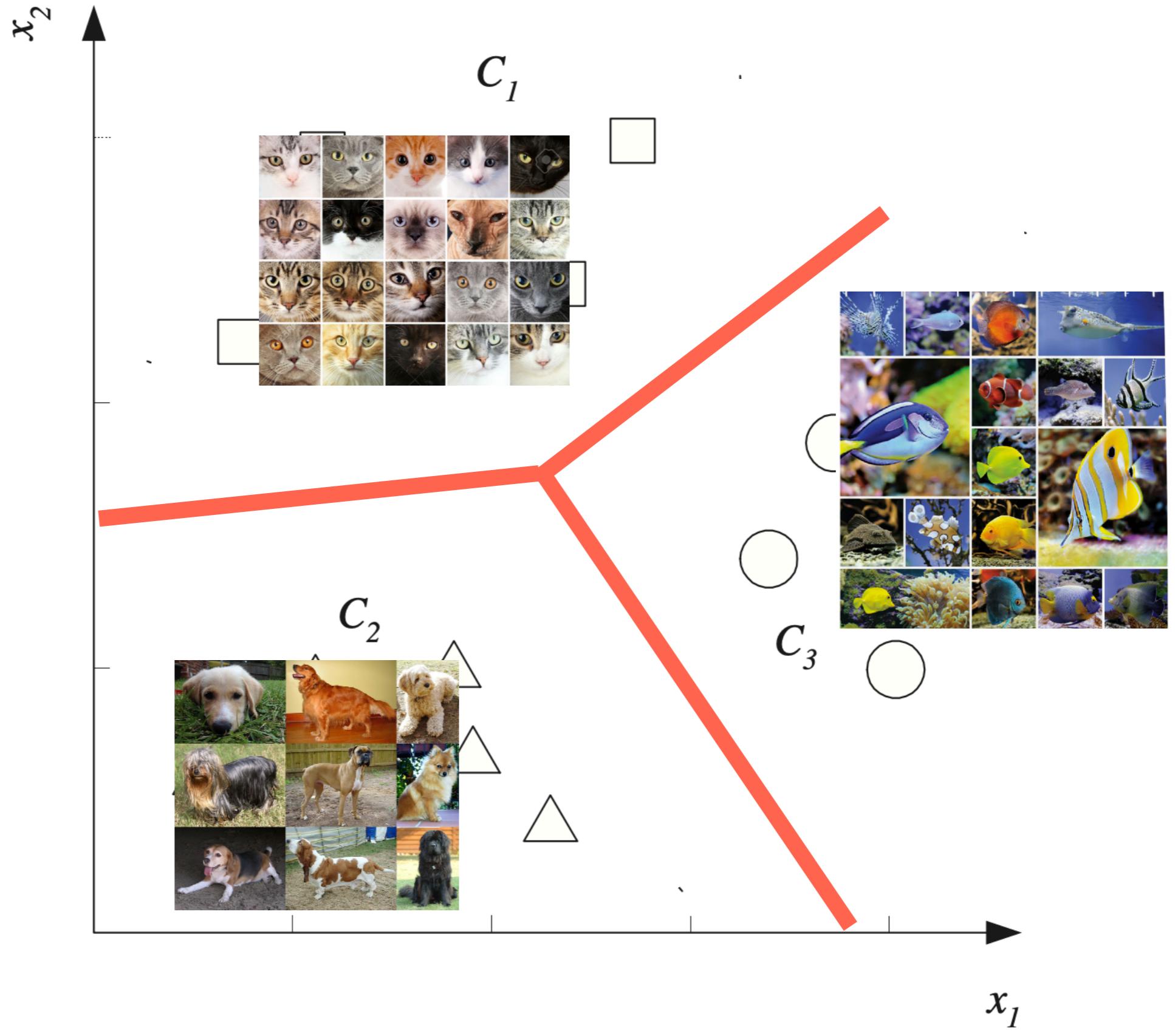
Classifier



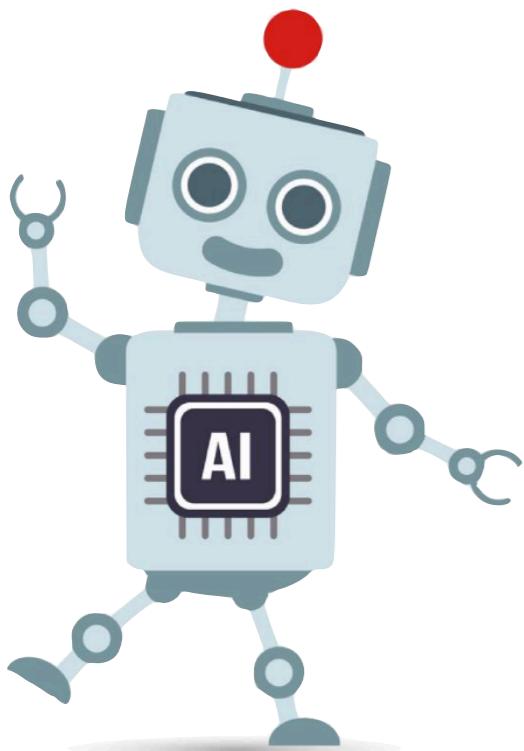
Classifier



Classifier

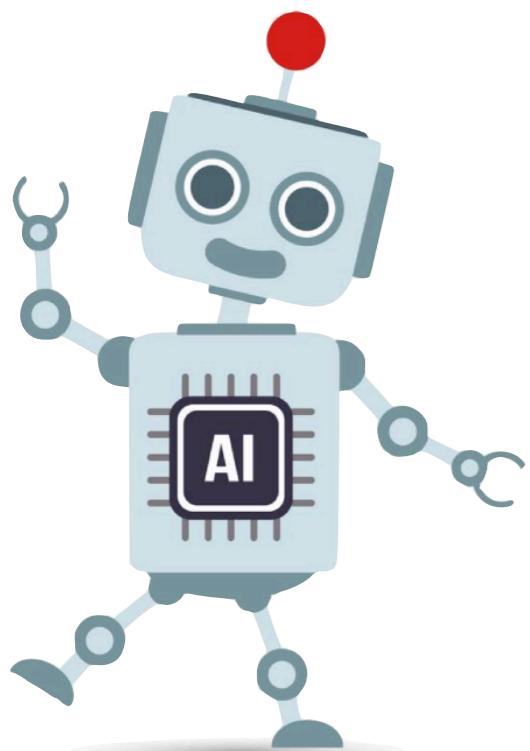


Goal: For input features \mathbf{x} ,
predict membership in
one of K classes.



Classifier

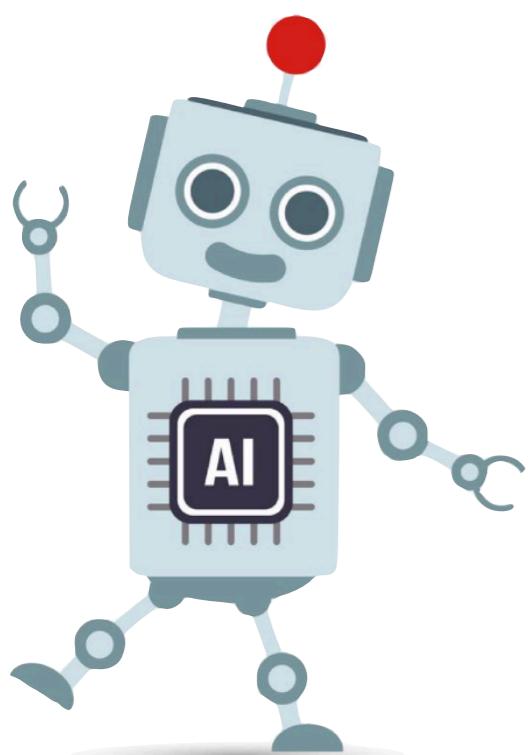
Goal: For input features x ,
predict membership in
one of K classes.



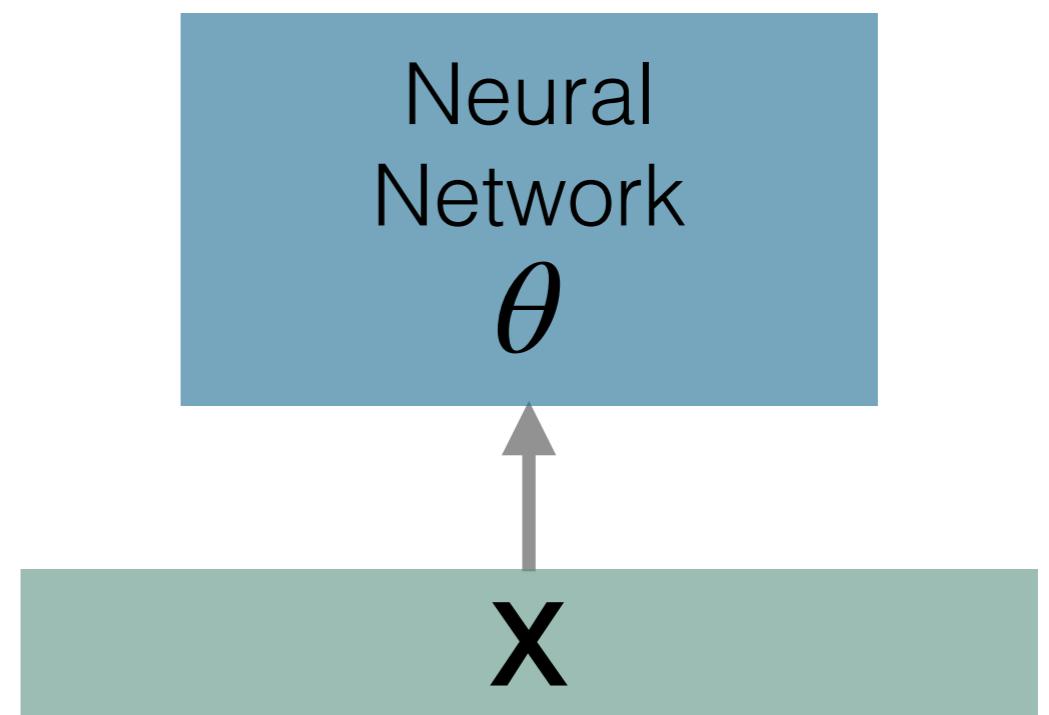
Classifier

x

Goal: For input features x ,
predict membership in
one of K classes.

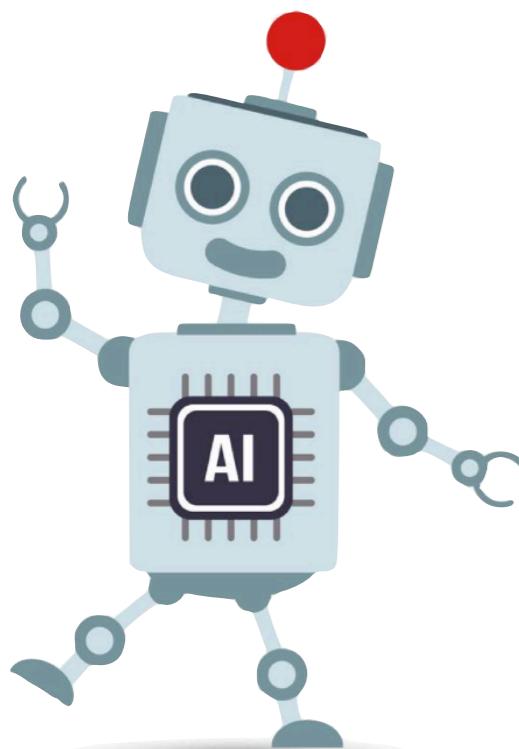


Classifier

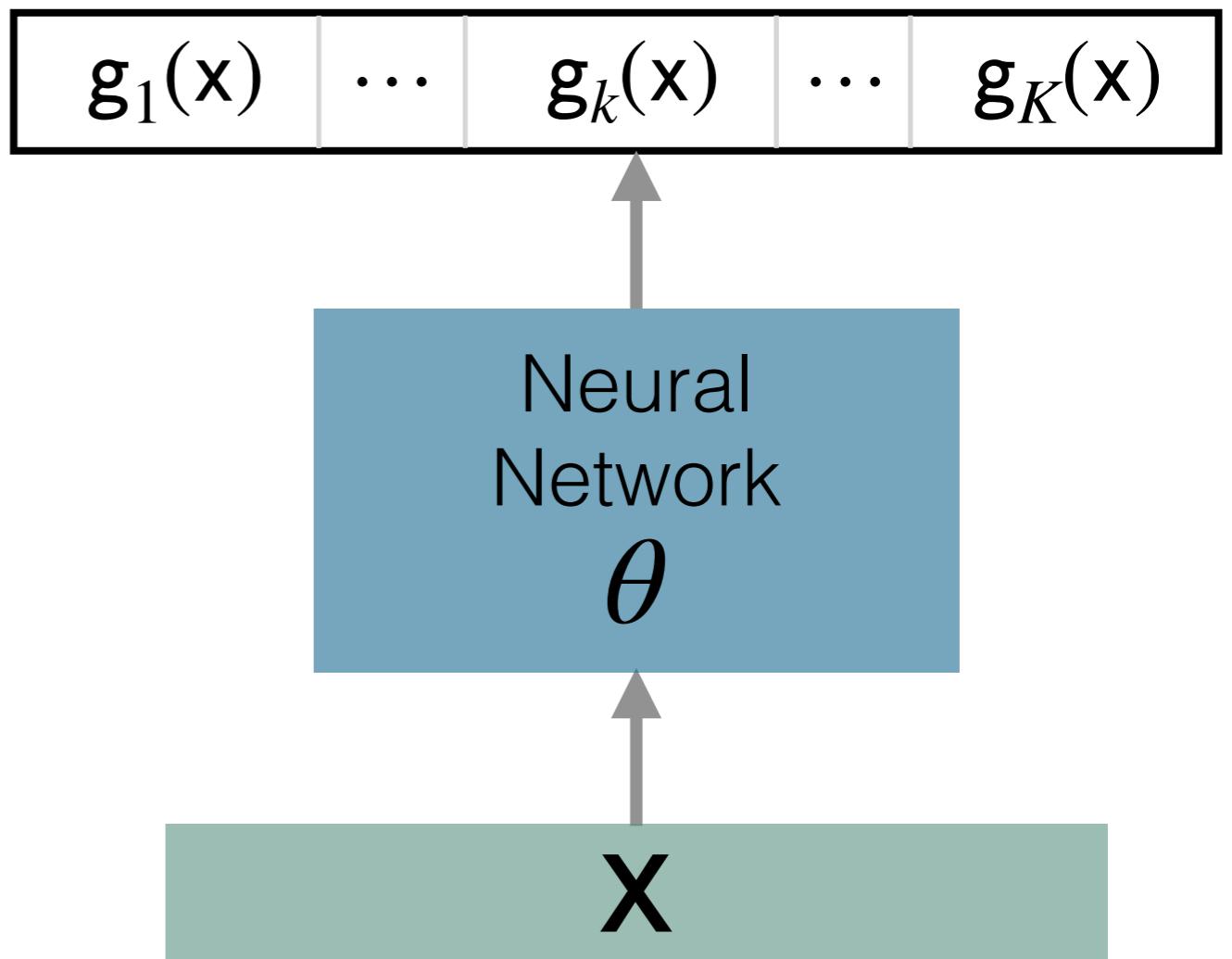


Goal: For input features x ,
predict membership in
one of K classes.

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

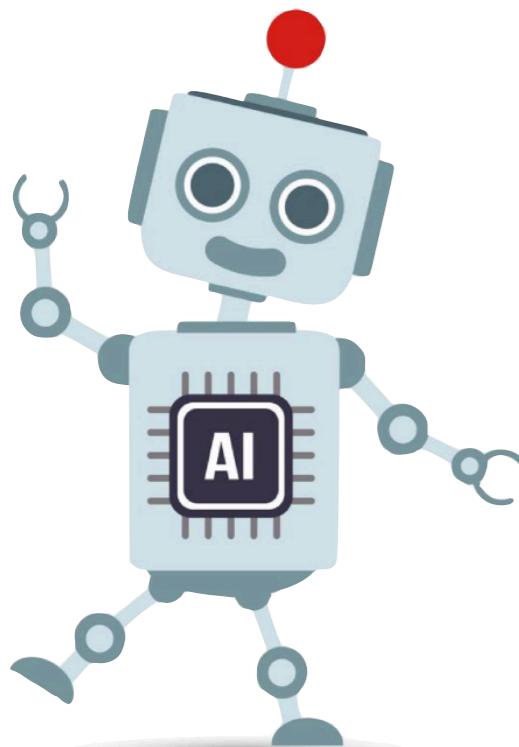


Classifier

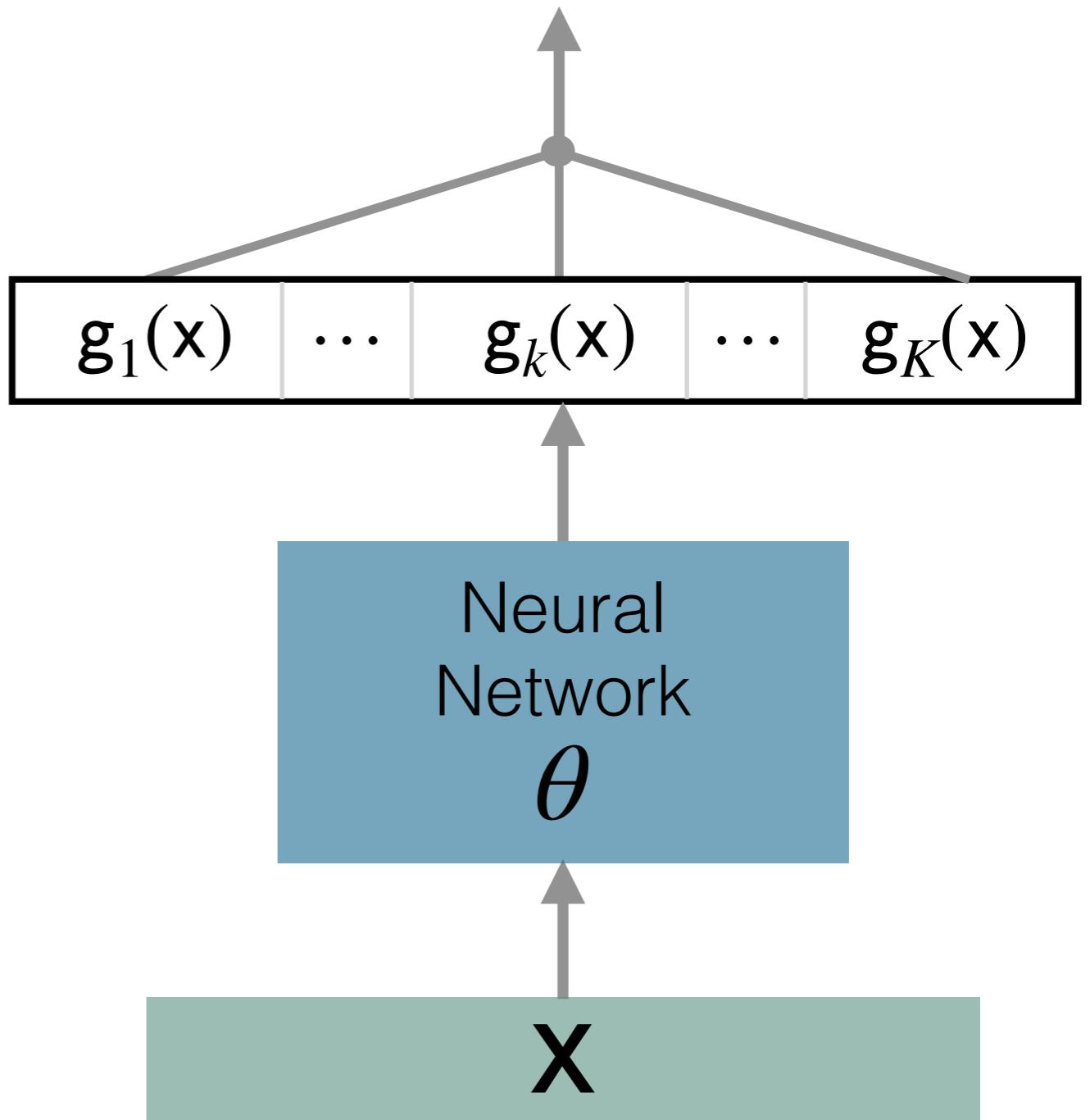


Goal: For input features x ,
predict membership in
one of K classes.

$$P(y|x) = \frac{\exp\{g_y(x)\}}{\sum_{k=1}^K \exp\{g_k(x)\}}$$

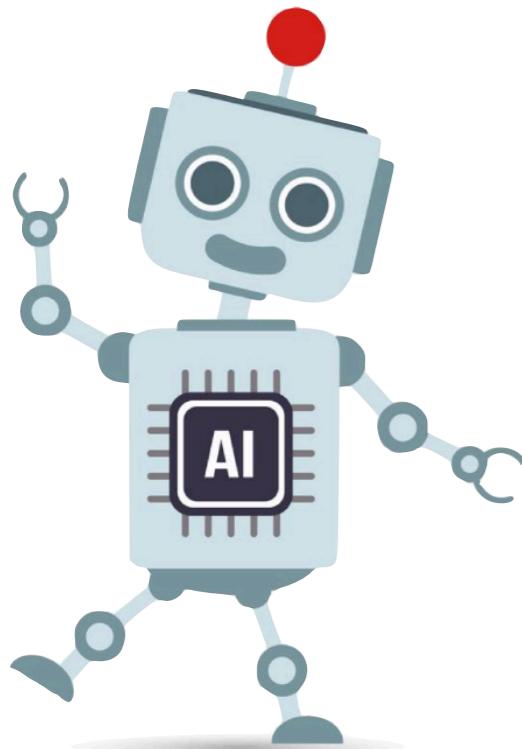


Classifier



Goal: For input features \mathbf{x} ,
predict membership in
one of K classes.

Training: Minimize with respect to θ ...

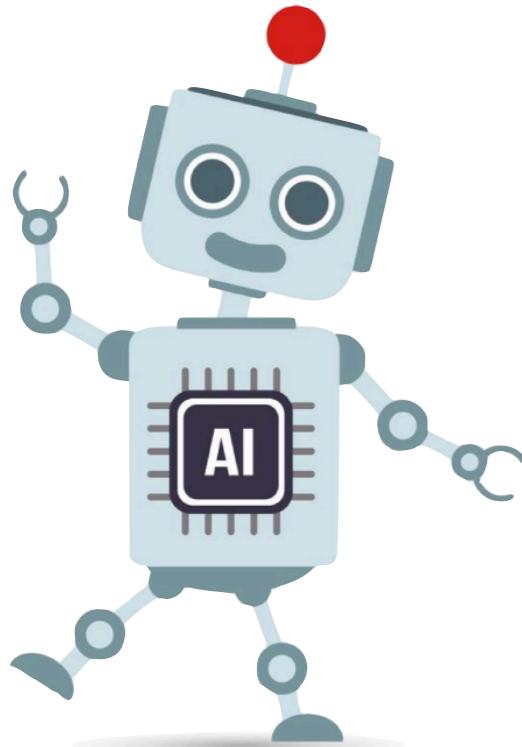


$$\ell(\theta; \mathbf{X}, \mathbf{y}) = -\log \left\{ \prod_n P(y_n | \mathbf{x}_n) \right\}$$
$$= -\sum_n \log \frac{\exp\{g_{y_n}(\mathbf{x}_n)\}}{\sum_{k=1}^K \exp\{g_k(\mathbf{x}_n)\}}$$

Classifier

Goal: For input features \mathbf{x} ,
predict membership in
one of K classes.

Training: Minimize with respect to θ ...

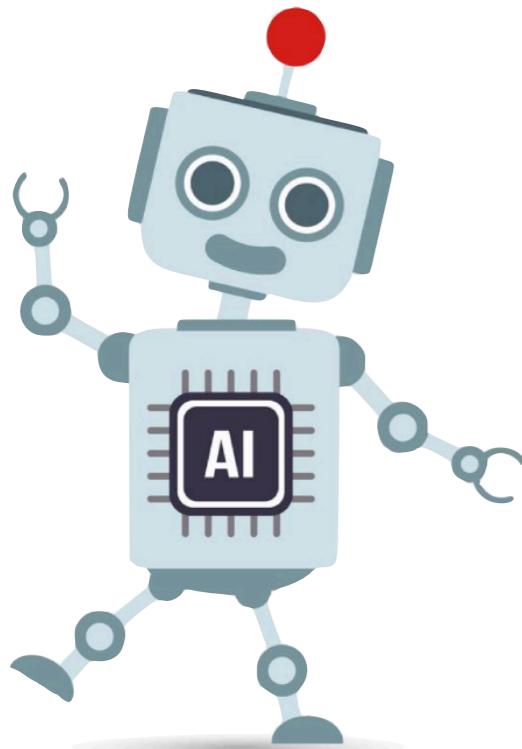


$$\ell(\theta; \mathbf{X}, \mathbf{y}) = -\log \left\{ \prod_n P(y_n | \mathbf{x}_n) \right\}$$
$$= -\sum_n \log \frac{\exp\{g_{y_n}(\mathbf{x}_n)\}}{\sum_{k=1}^K \exp\{g_k(\mathbf{x}_n)\}}$$

Classifier

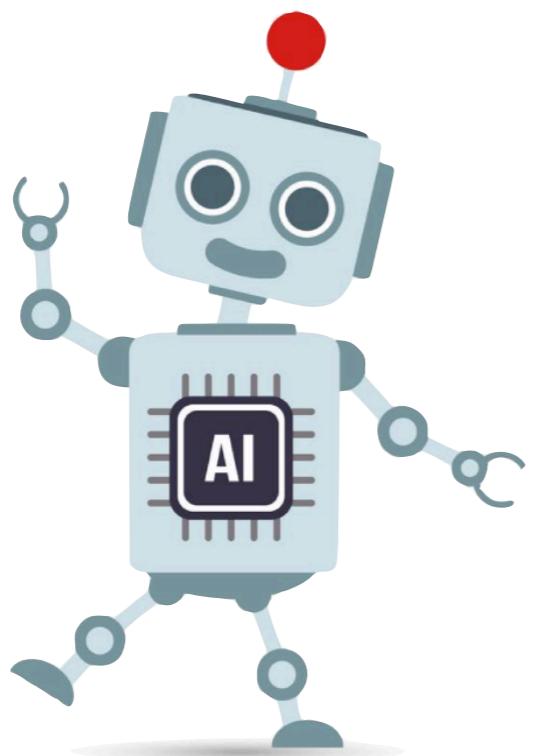
Goal: For input features \mathbf{x} ,
predict membership in
one of K classes.

Training: Minimize with respect to θ ...

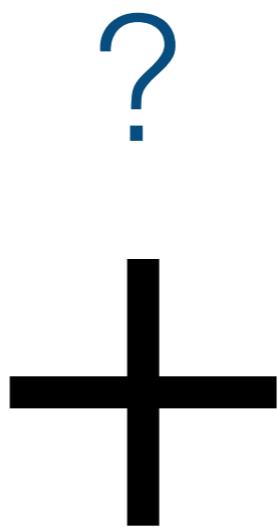


$$\ell(\theta; \mathbf{X}, \mathbf{y}) = -\log \left\{ \prod_n P(y_n | \mathbf{x}_n) \right\}$$
$$= - \sum_n \log \frac{\exp\{g_{y_n}(\mathbf{x}_n)\}}{\sum_{k=1}^K \exp\{g_k(\mathbf{x}_n)\}}$$

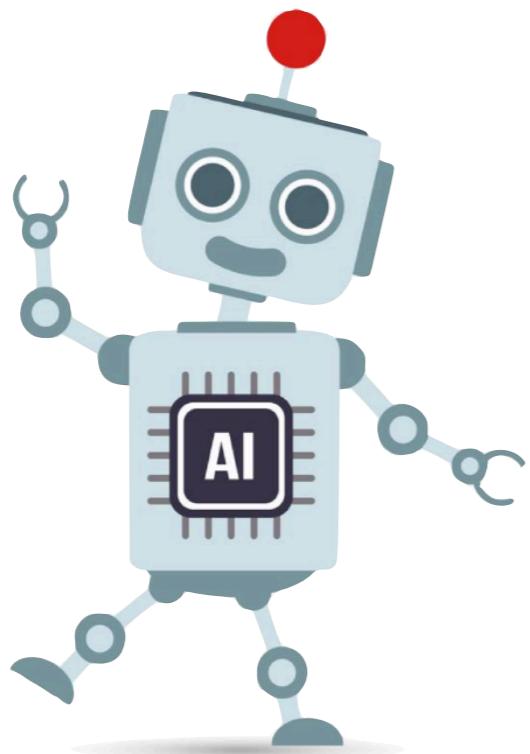
Classifier



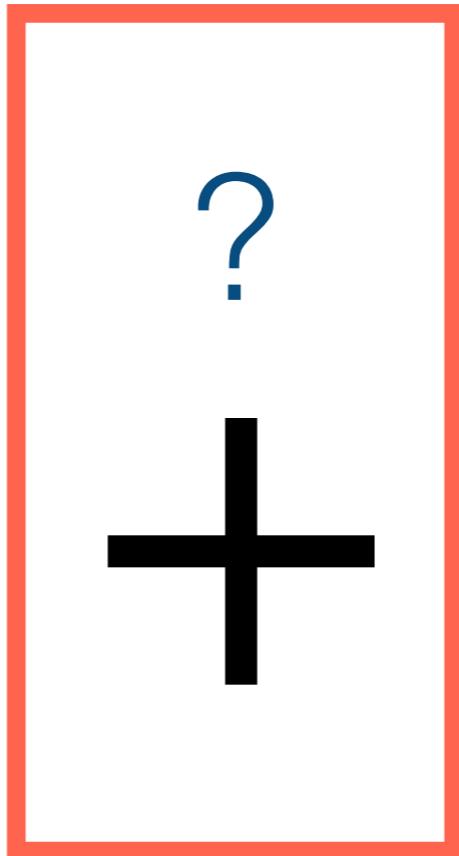
Classifier



Human Expert



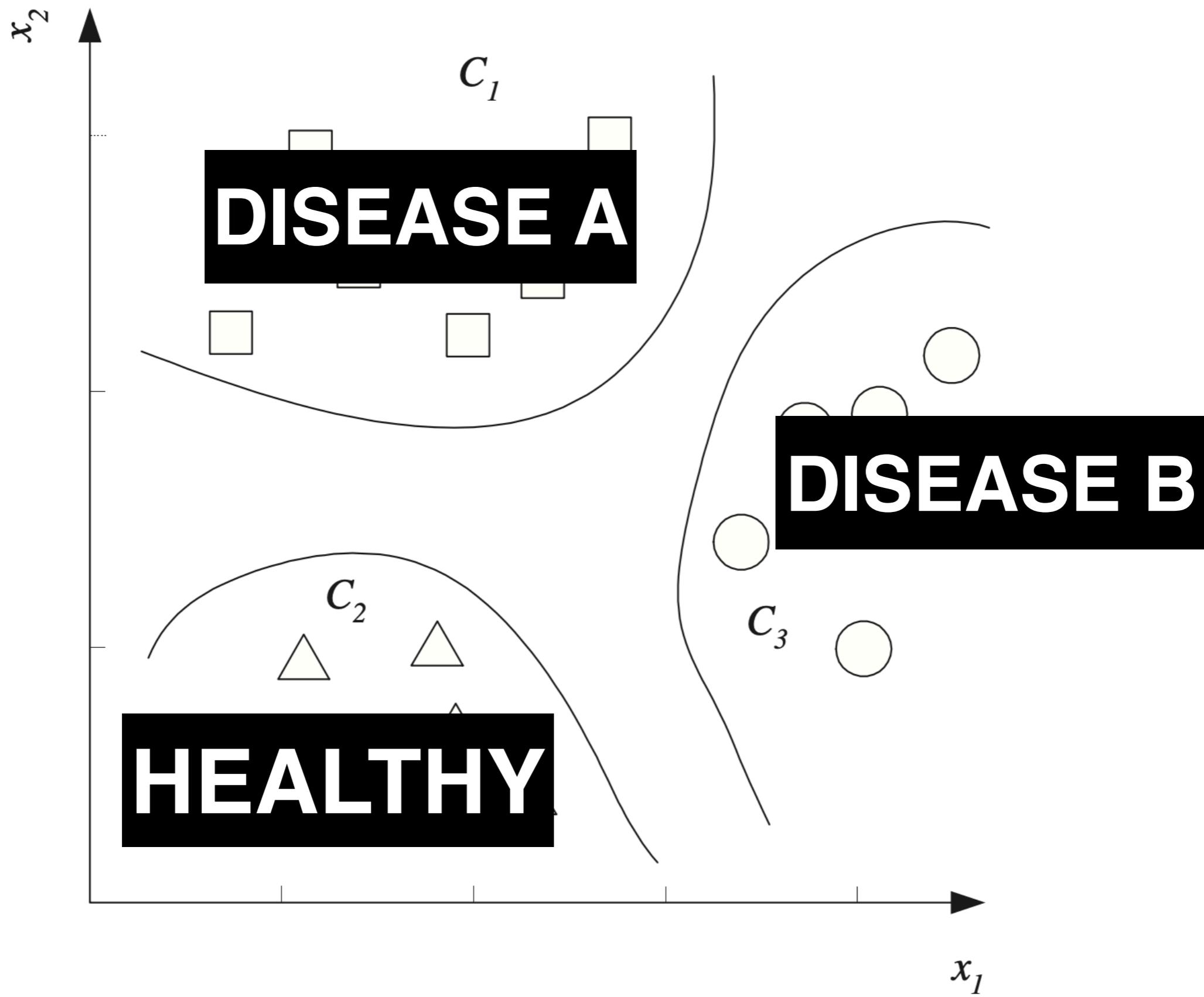
Classifier

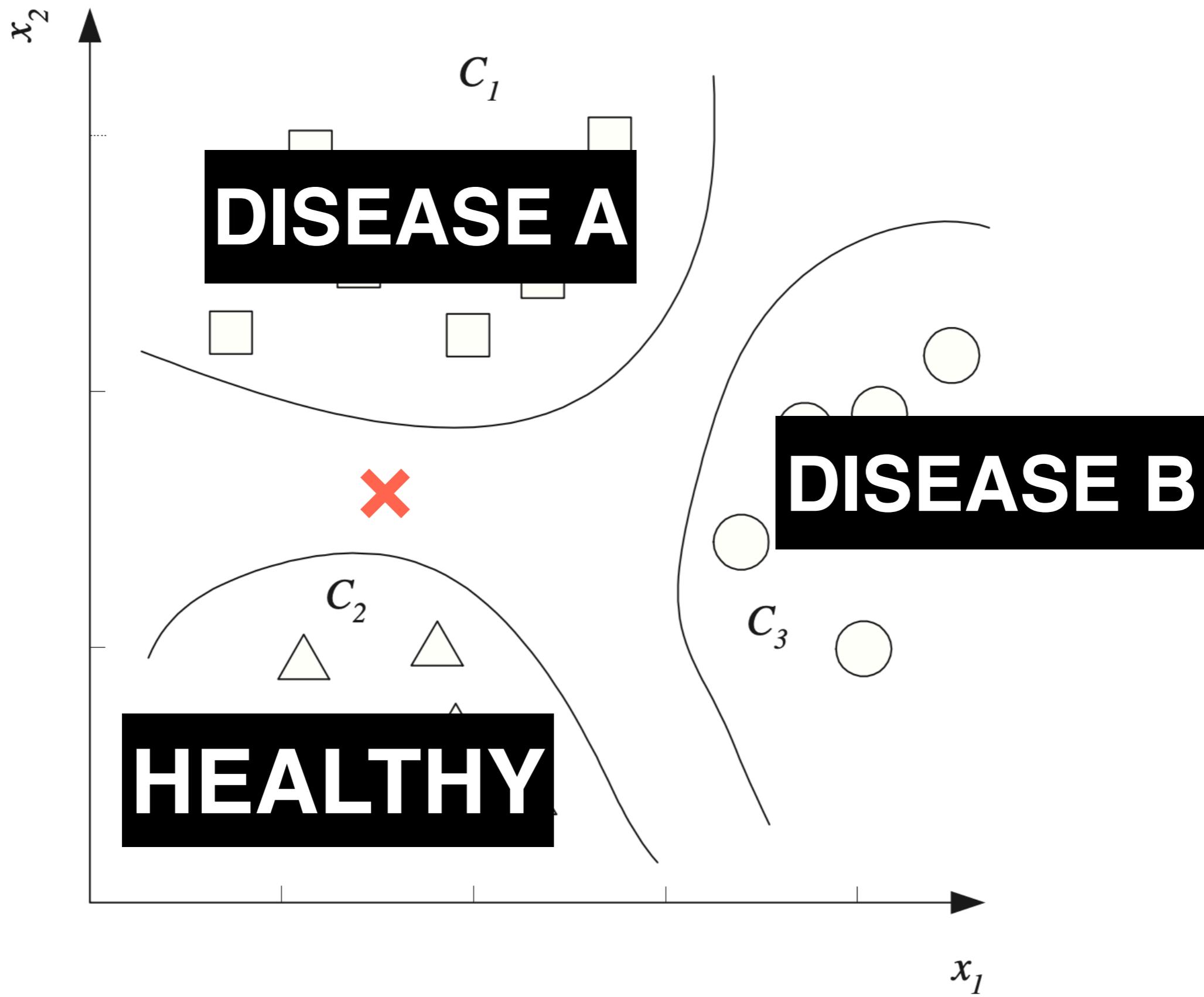


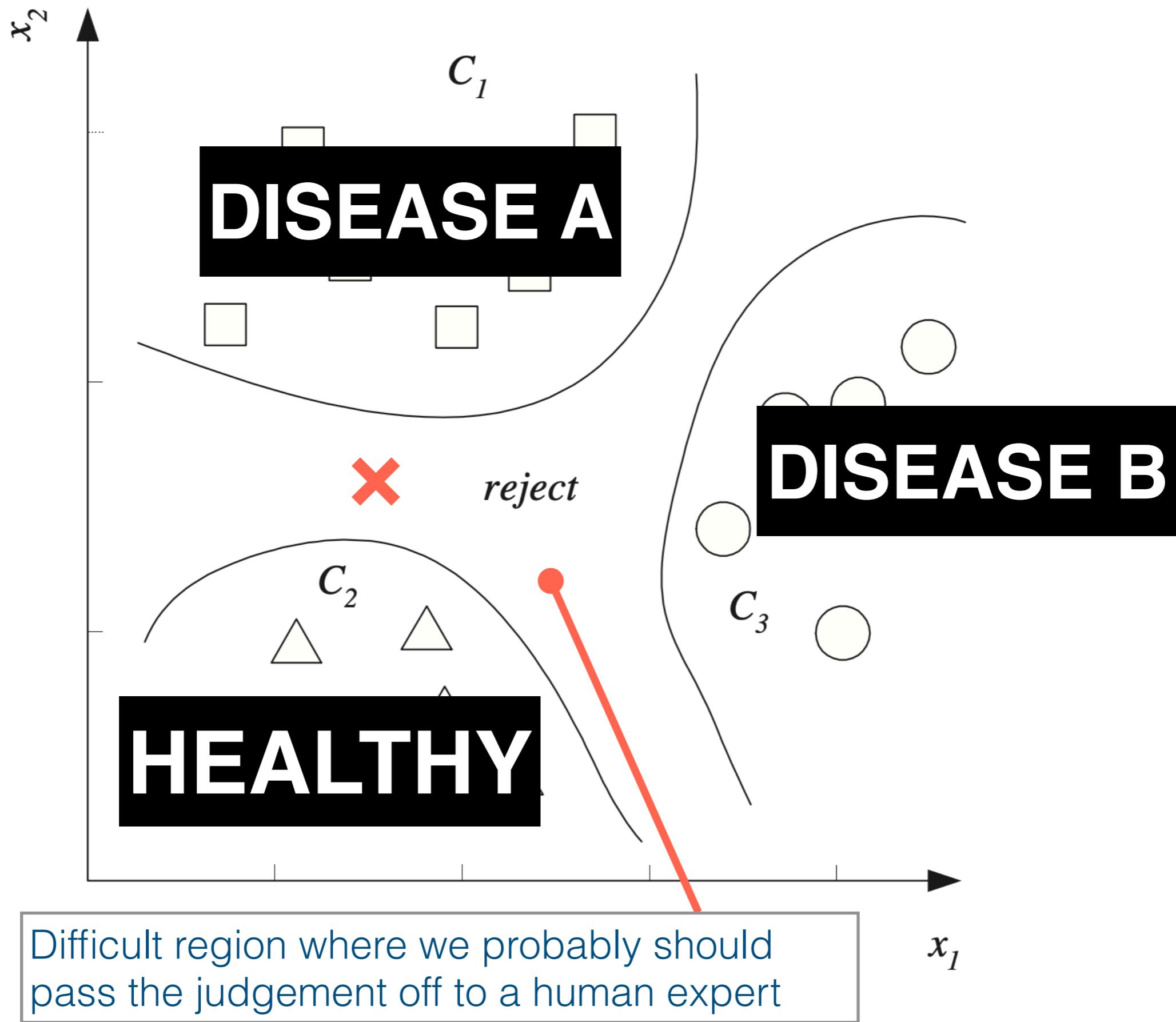
Human Expert

Warm Up

Classification with a
Rejection Option

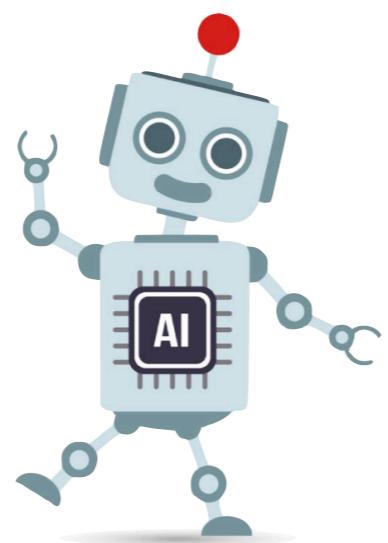








X

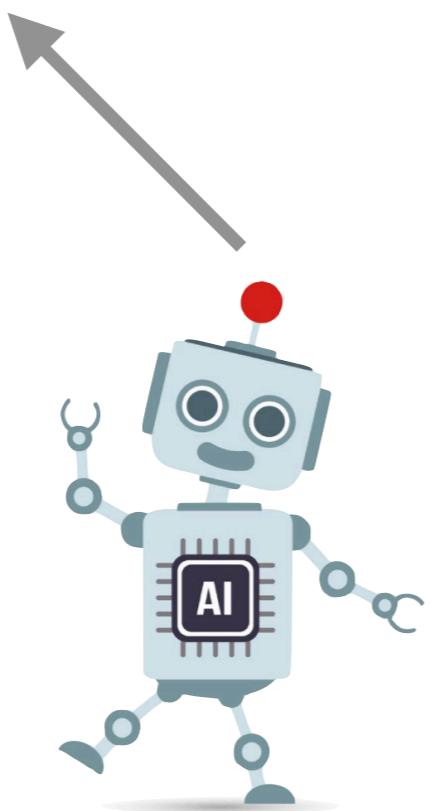


Classifier

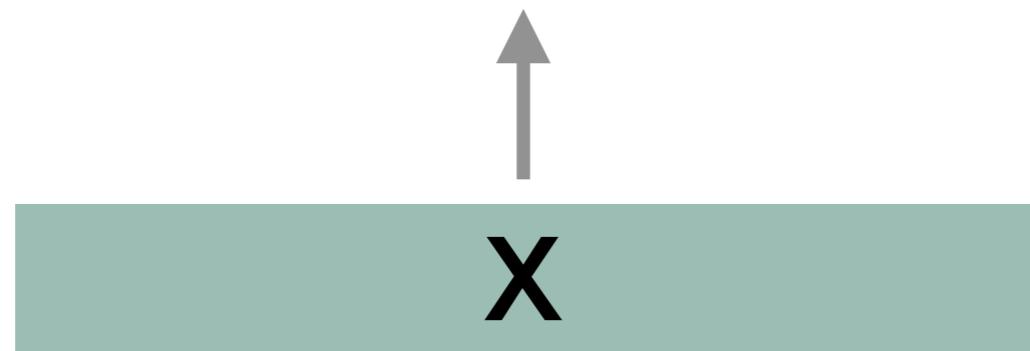


X

make
prediction

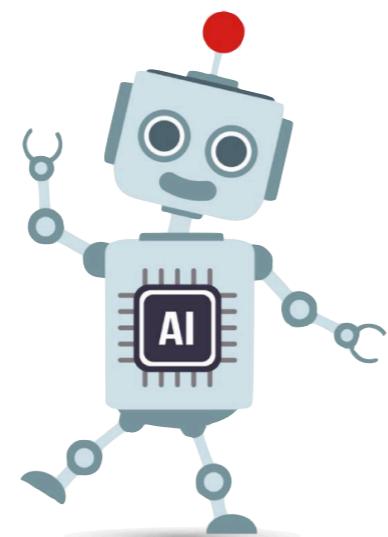


Classifier



make
prediction

abstain

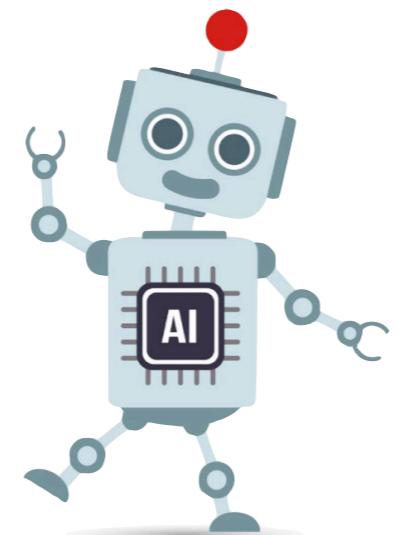


Classifier

X

make
prediction

abstain

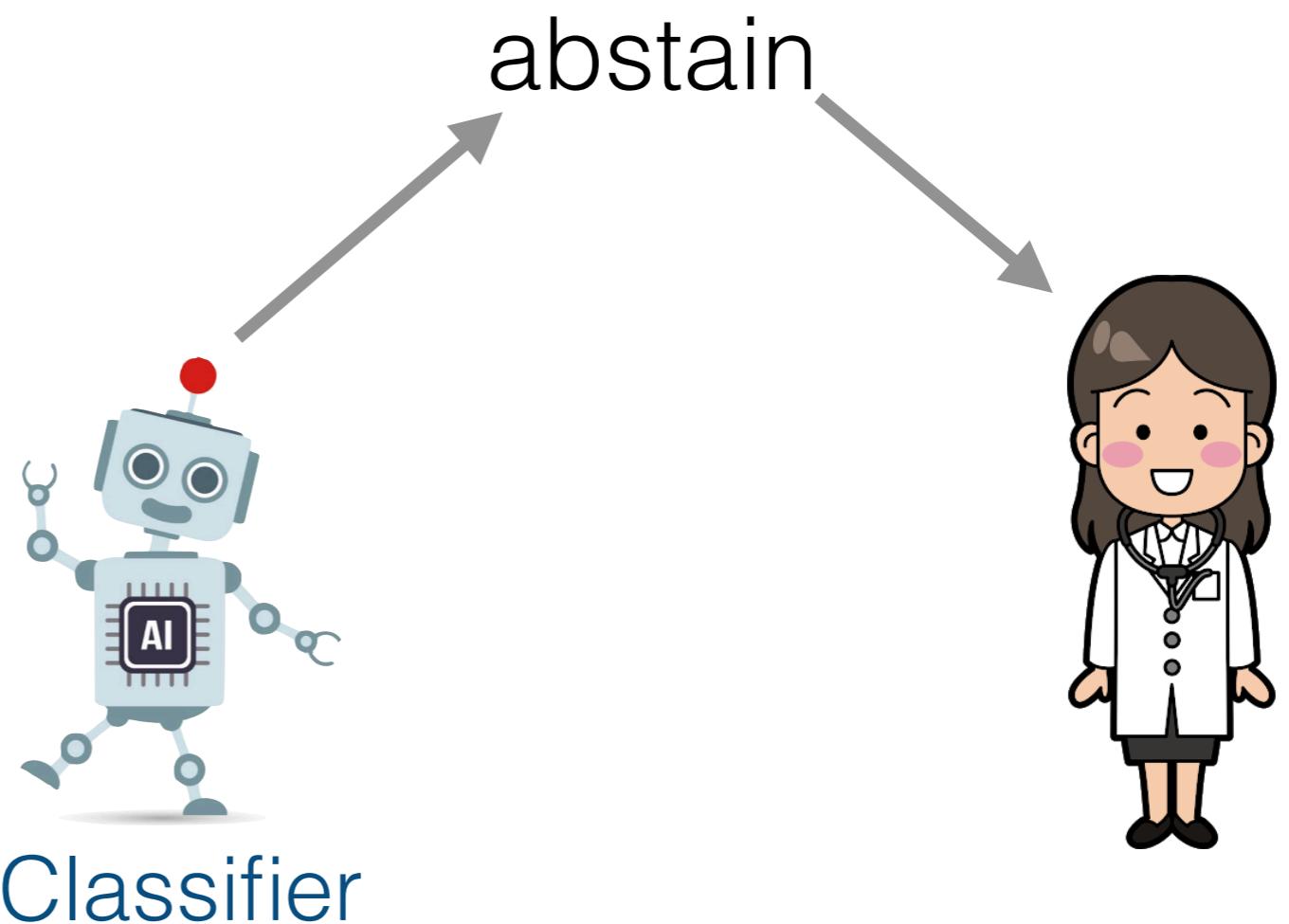


Classifier

X

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

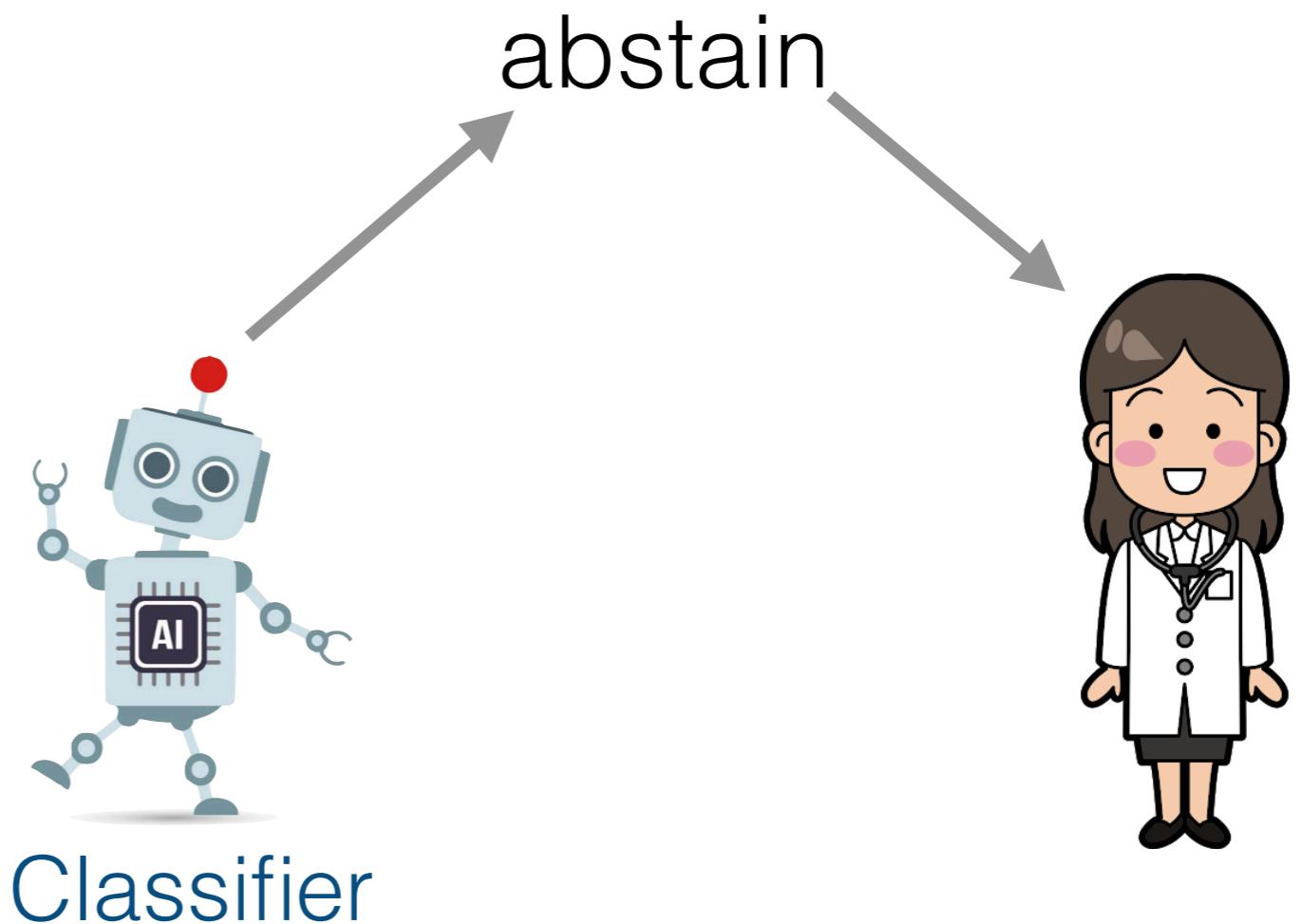
$$\max_y P(y | x_n) < \tau$$



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

$$\max_y P(y | x_n) < \tau$$

Human behavior
is not modeled!



Challenge: how can we model the human?

If they are a true expert, modeling their decision making— $P_h(y | x)$ —is assumed to be impossible.

Learning to Defer to an Expert

[Madras et al., NeurIPS 2018]

Better Formulation

*Model what the human knows,
so we can enable collaboration*

Better Formulation

*Model what the human knows,
so we can enable collaboration*

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

↑
expert predictions

Better Formulation

*Model what the human knows,
so we can enable collaboration*

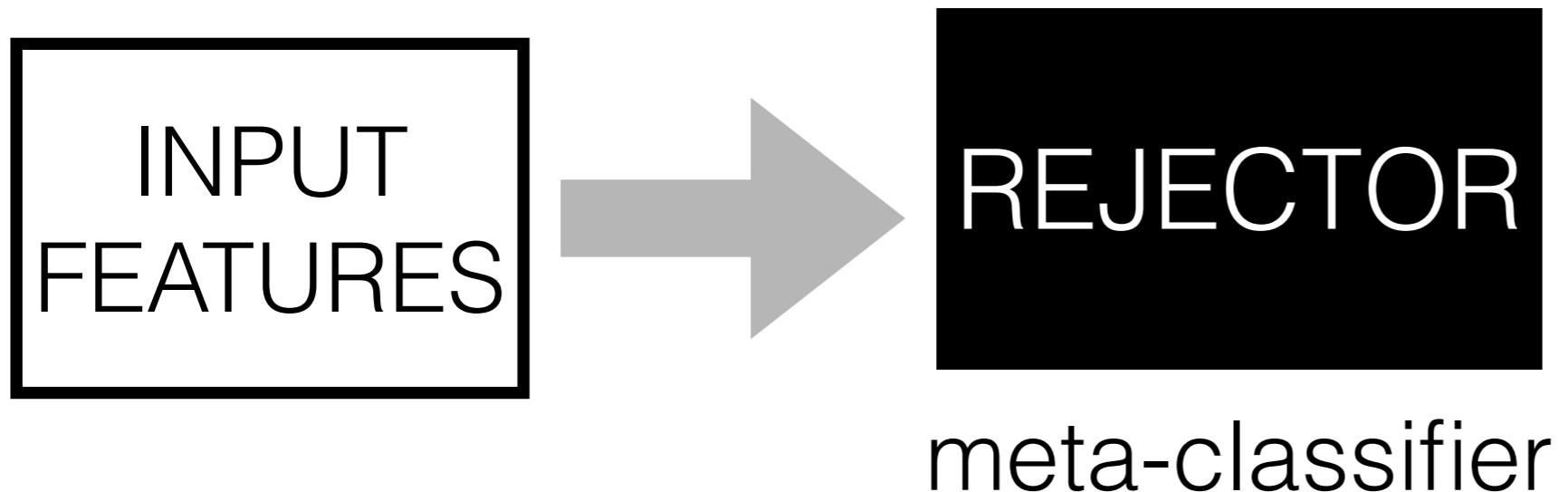
Data: $\mathcal{D} = \{\mathbf{x}_n, \mathbf{y}_n, \mathbf{m}_n\}_{n=1}^N$


expert predictions

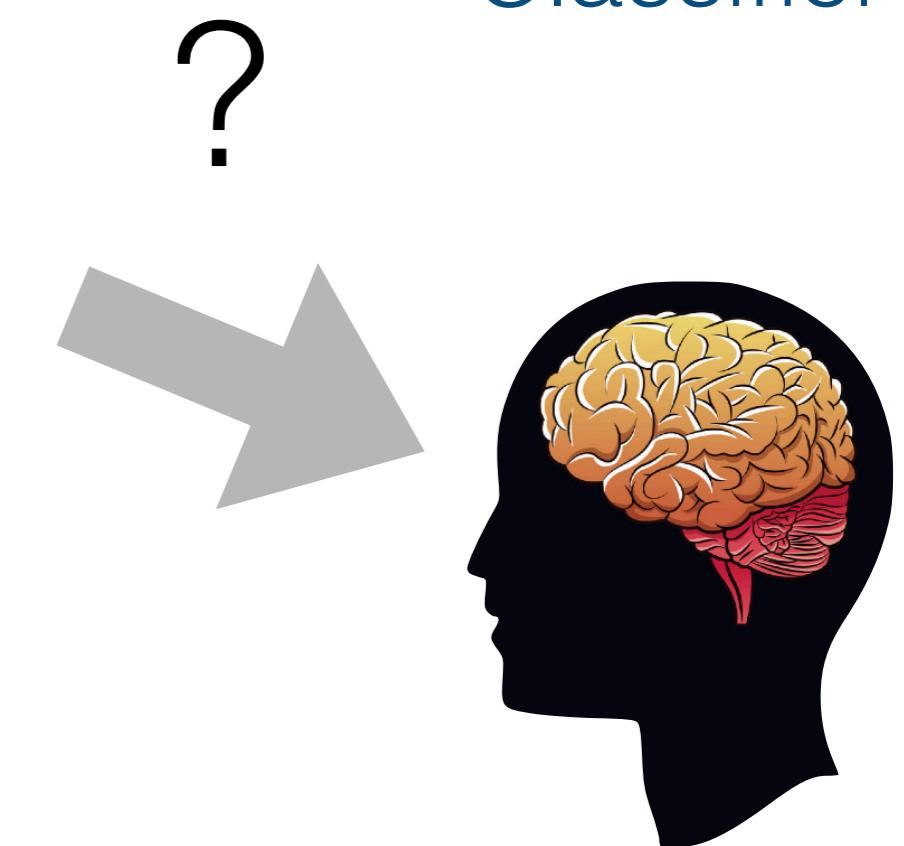
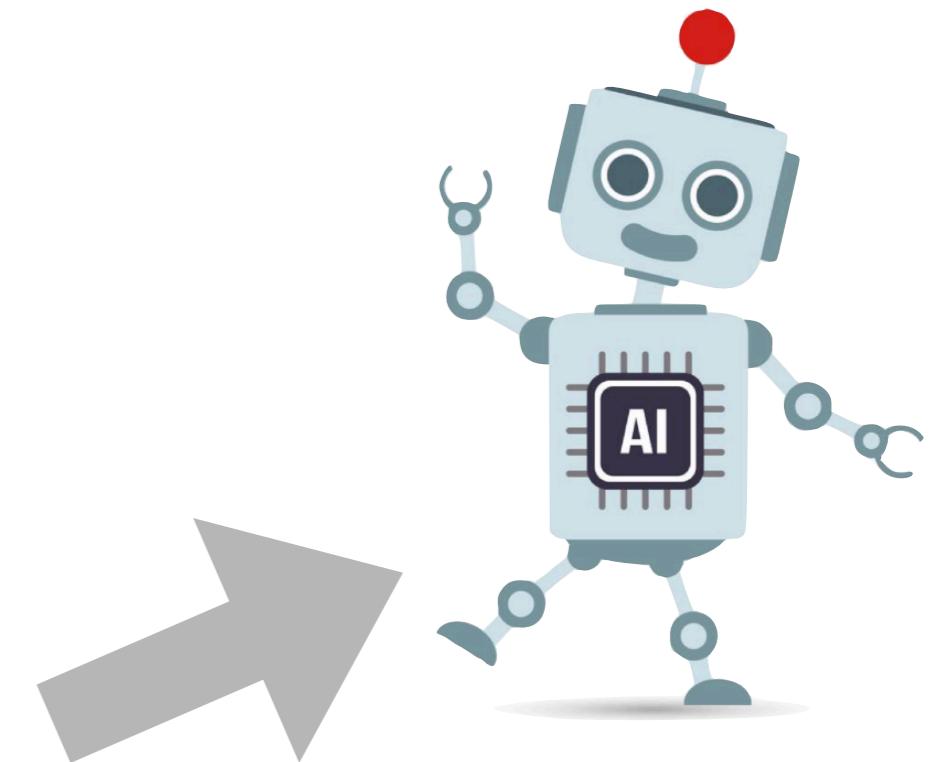
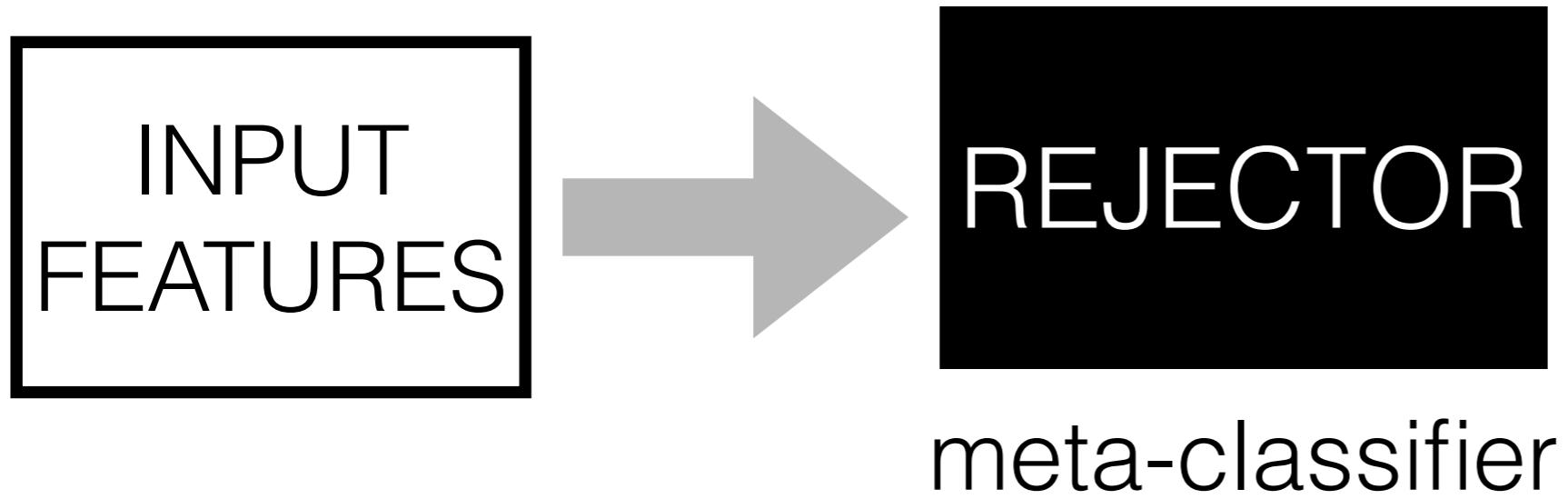
Models:

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

Learning to Defer

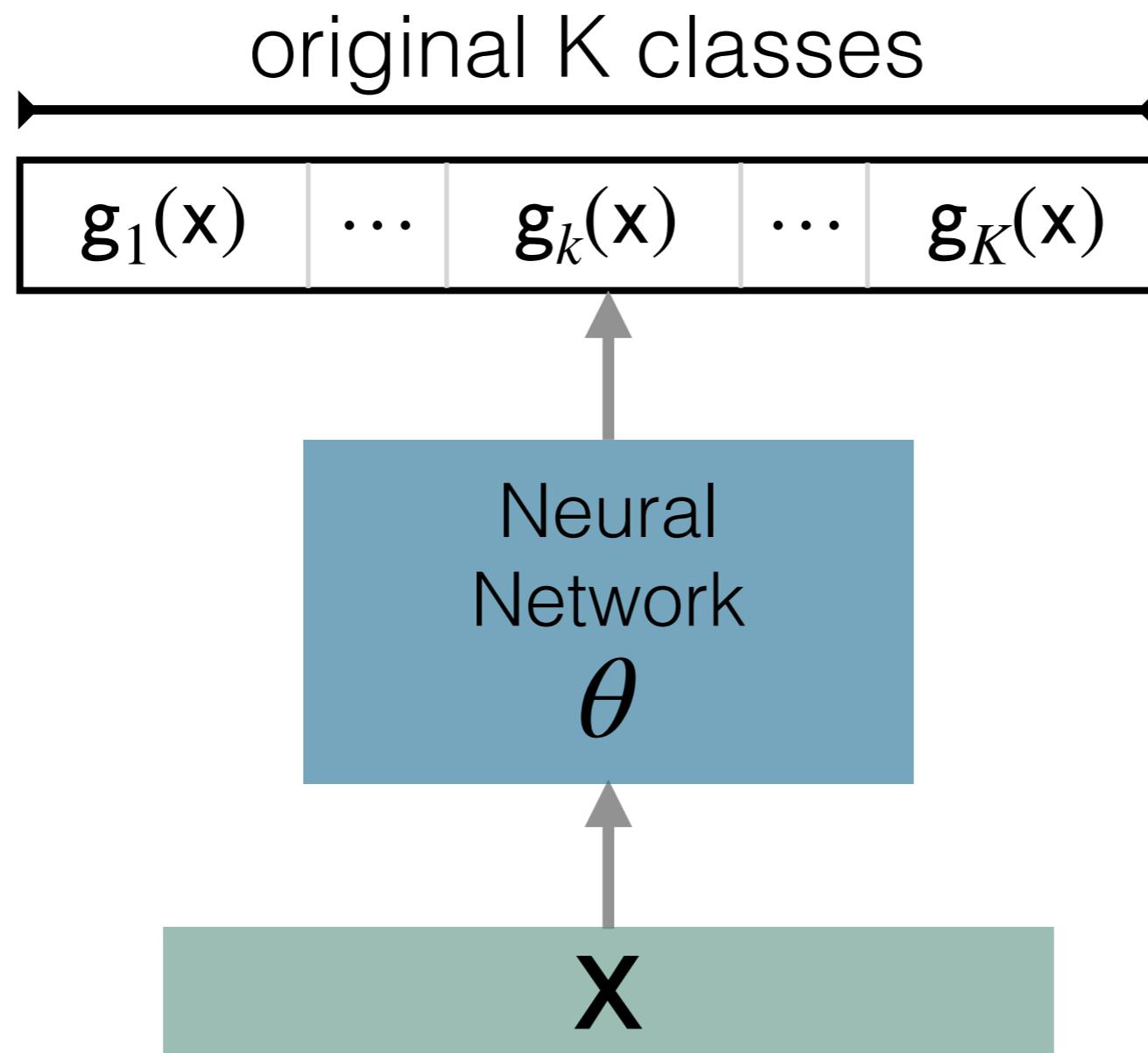


Learning to Defer

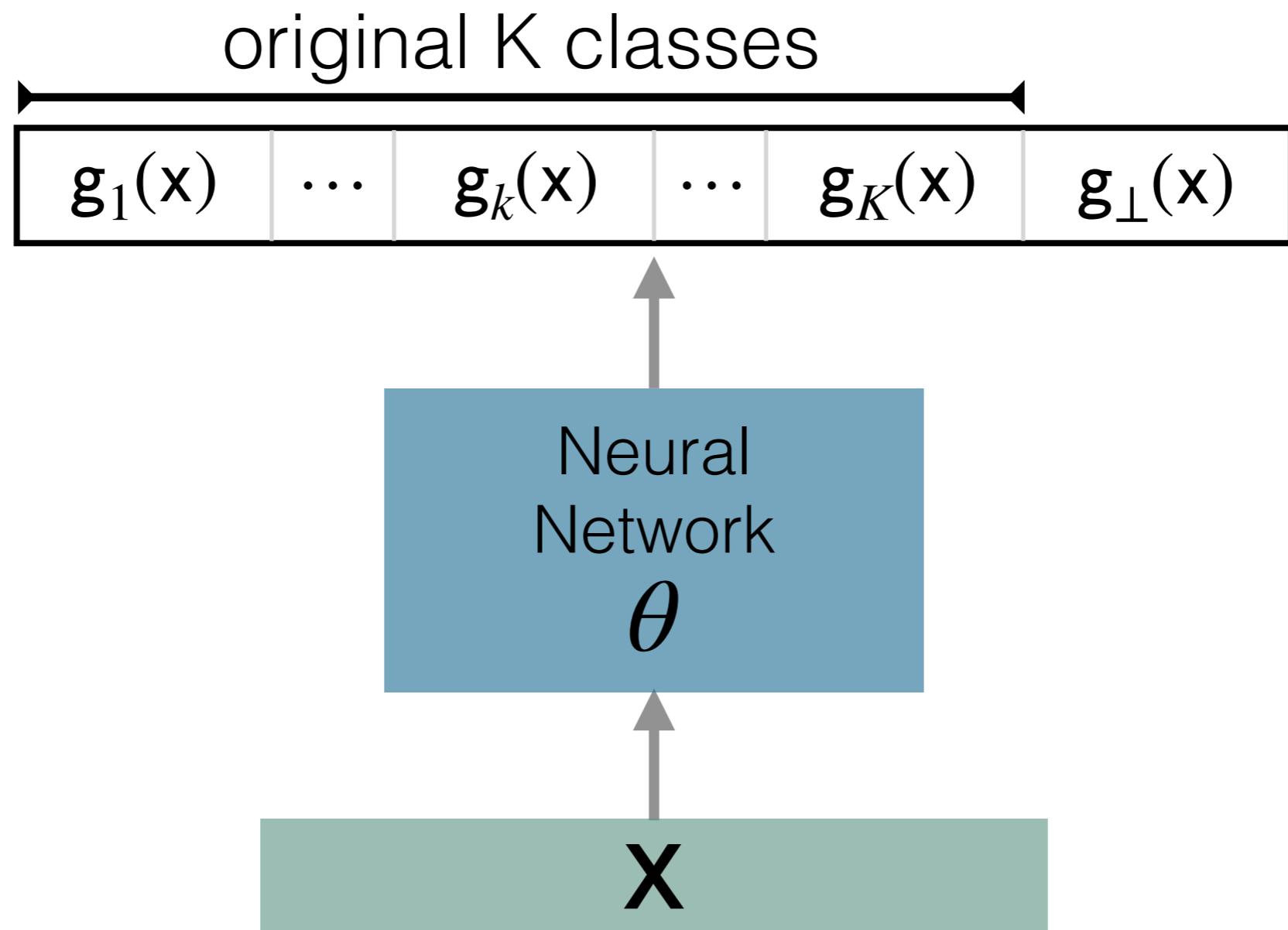


Softmax Approach [Mozannar & Sontag, ICML 2020]

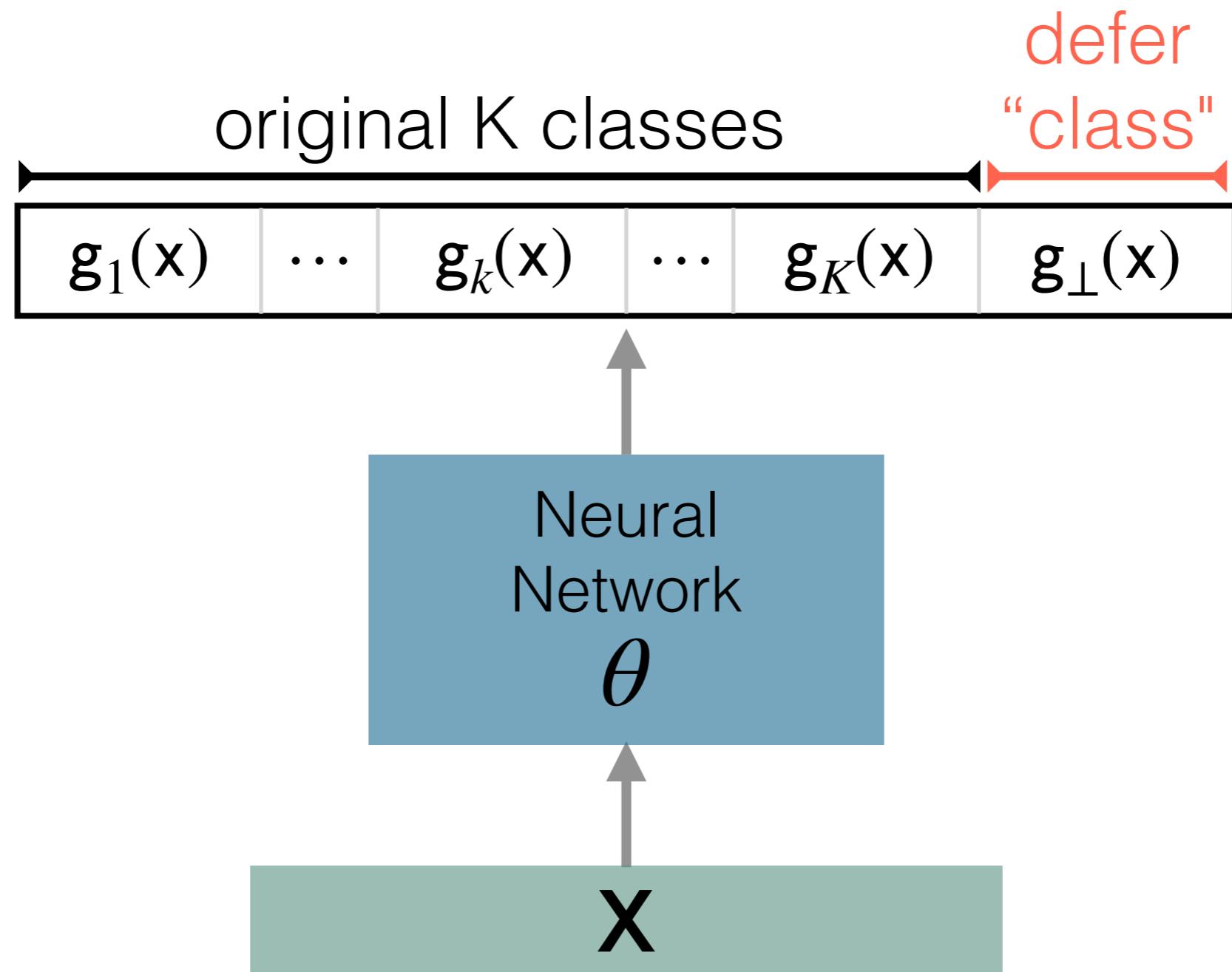
Softmax Approach [Mozannar & Sontag, ICML 2020]



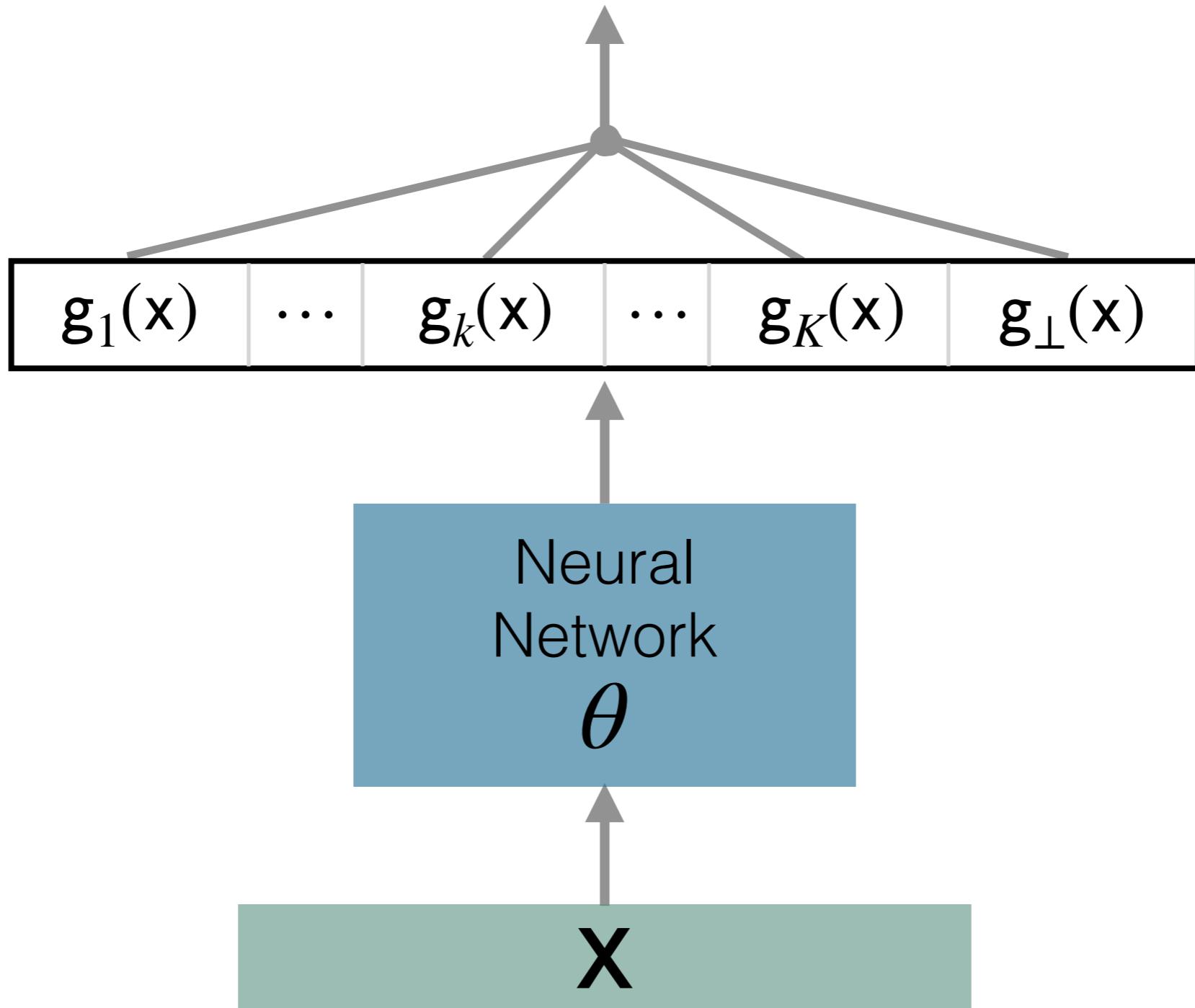
Softmax Approach [Mozannar & Sontag, ICML 2020]



Softmax Approach [Mozannar & Sontag, ICML 2020]

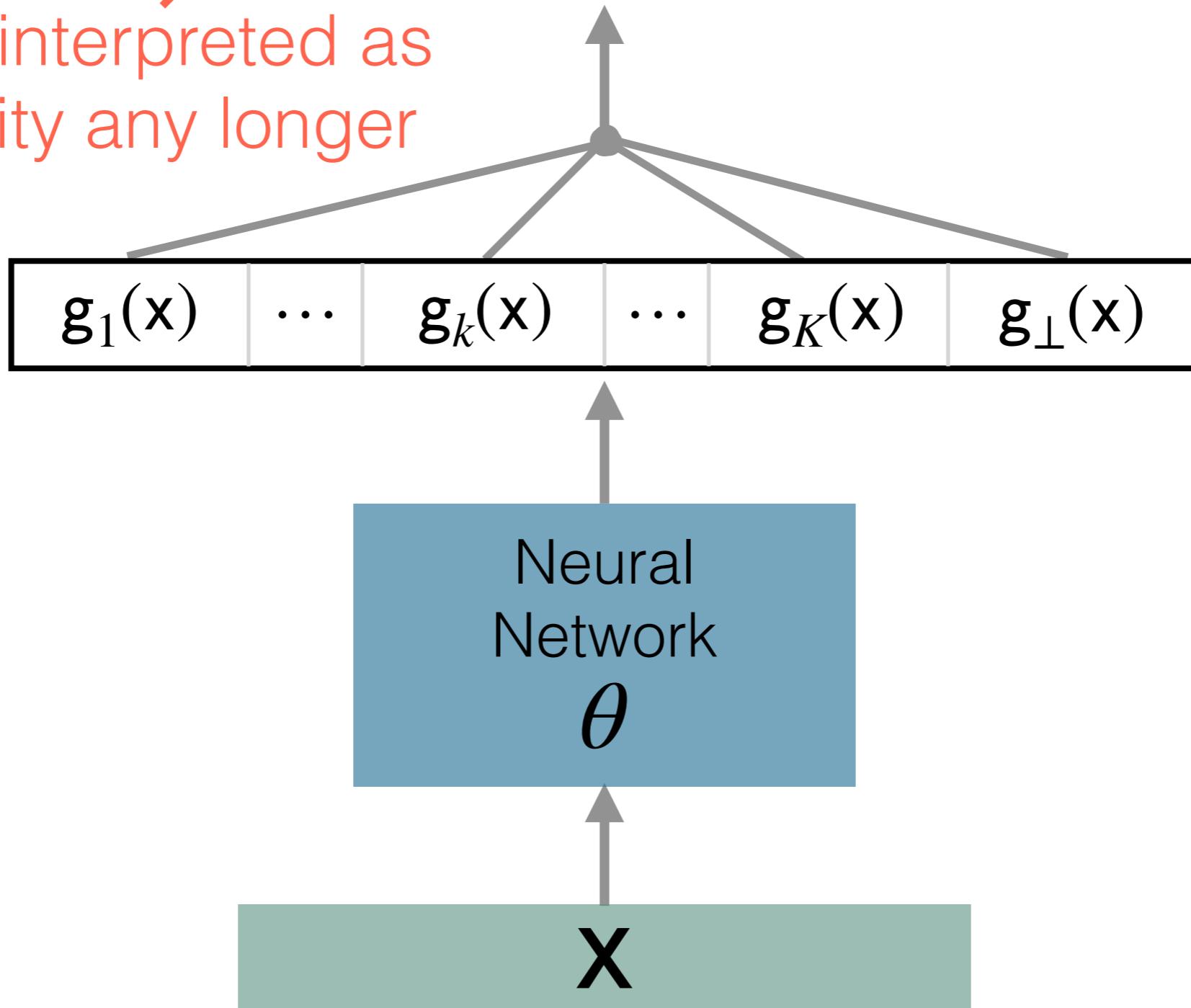


$$p_i(x) = \frac{\exp\{g_i(x)\}}{\sum_{k=1}^{K+1} \exp\{g_k(x)\}}$$

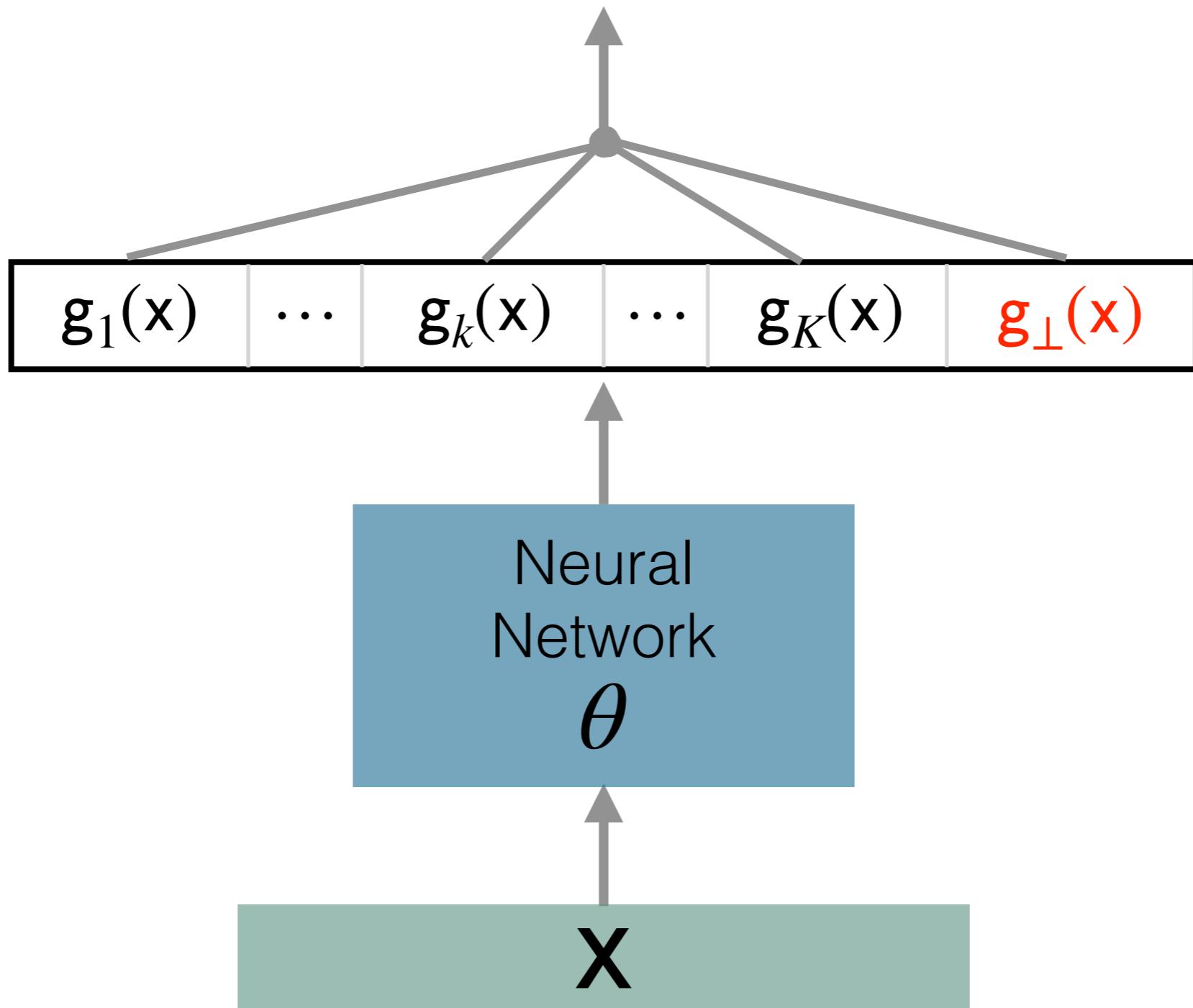


$$p_i(x) = \frac{\exp\{g_i(x)\}}{\sum_{k=1}^{K+1} \exp\{g_k(x)\}}$$

cannot be interpreted as
a probability any longer



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



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

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

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

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

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

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.

Consistency

Classifier and rejector have the following Bayes optimal solutions (under 0-1 loss):

Consistency

Classifier and rejector have the following Bayes optimal solutions (under 0-1 loss):

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

Consistency

Classifier and rejector have the following Bayes optimal solutions (under 0-1 loss):

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

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

probability that the expert is correct

Our work

Is this learning-to-defer
system calibrated?

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?



Human-computer collaboration for skin cancer recognition

Philipp Tschandl^{1,17}, Christoph Rinner^{1,17}, Zoe Apalla³, Giuseppe Argenziano^{1,4}, Noel Codella⁵, Allan Halpern⁶, Monika Janda⁷, Aimilios Lallas³, Caterina Longo^{8,9}, Josep Malvehy^{10,11}, John Paoli^{12,13}, Susana Puig^{10,11}, Cliff Rosendahl¹⁴, H. Peter Soyer¹⁵, Iris Zalaudek¹⁶ and Harald Kittler¹✉

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^{10–15}. 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.”

Recall...

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

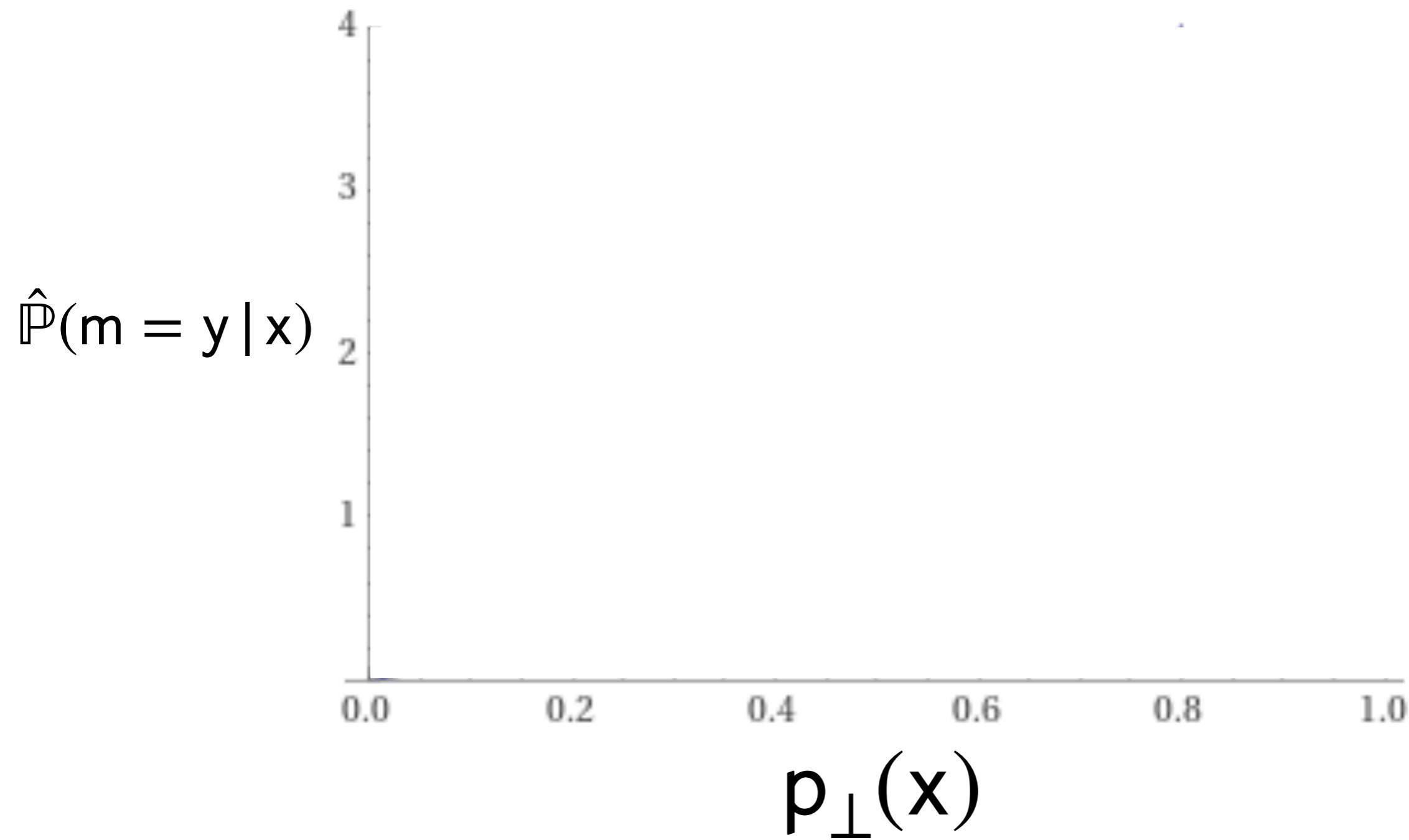
Recall...

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

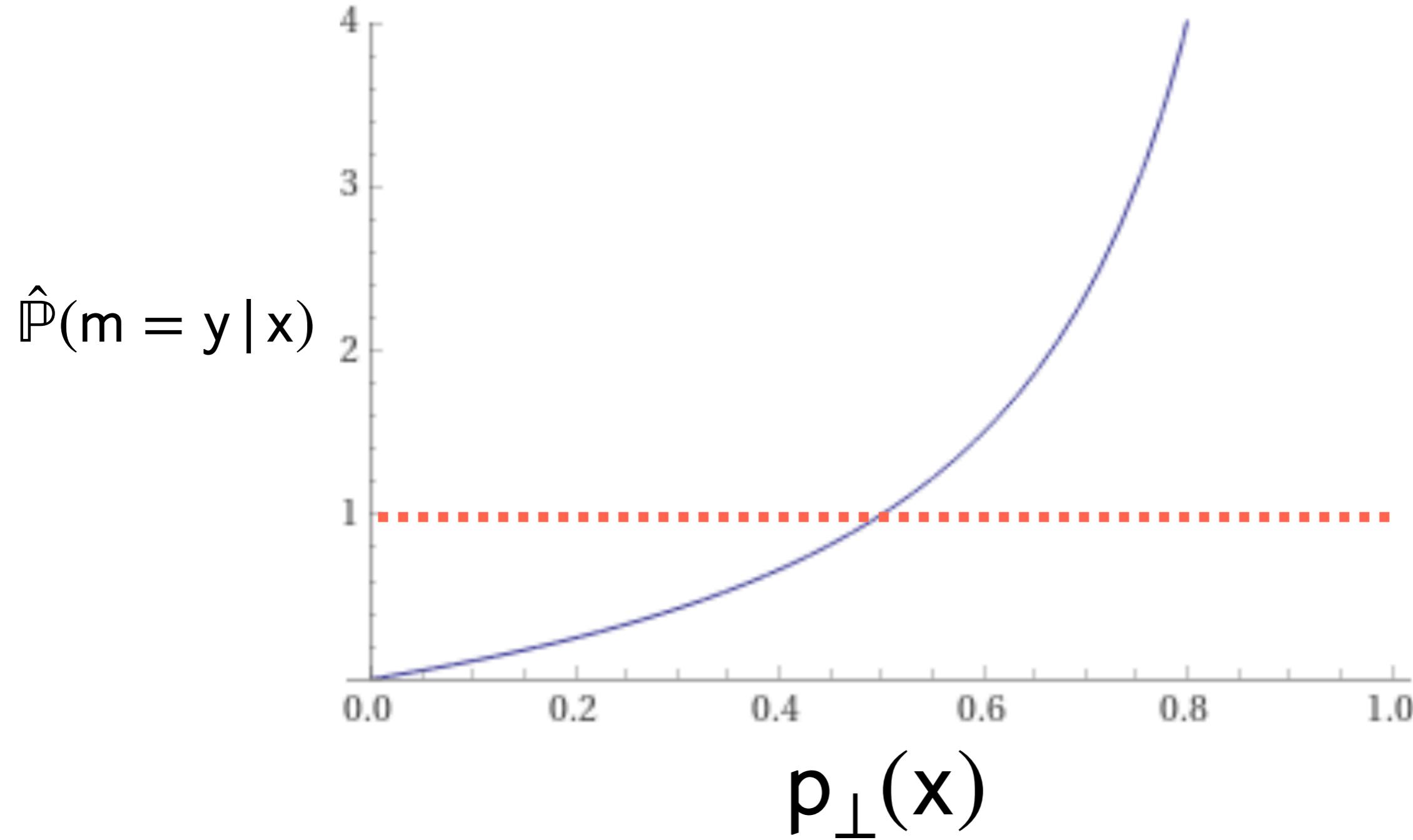
From Mozannar's & Sontag's Theorem #1...

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

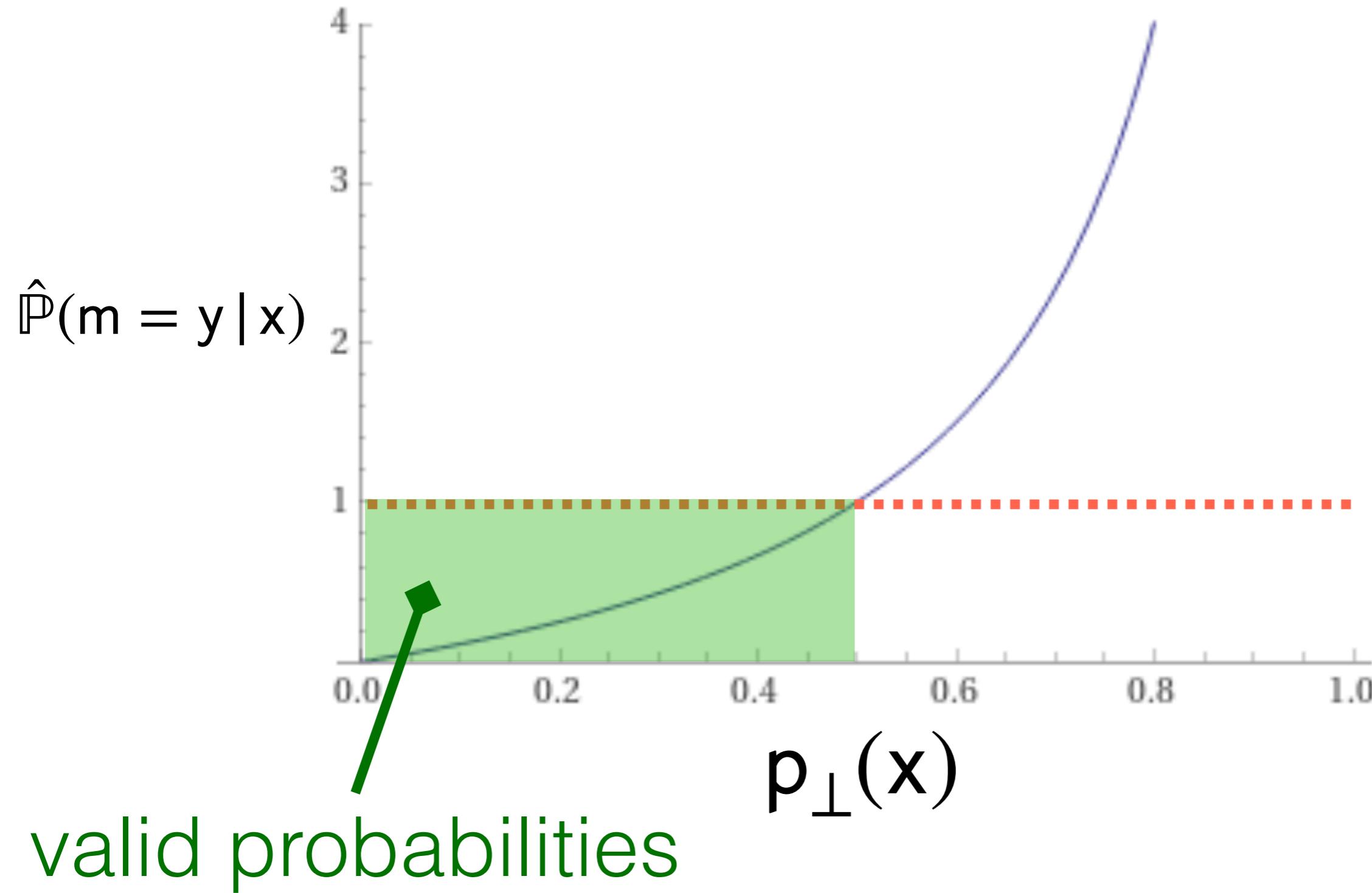
Estimating $P(m = y | x)$



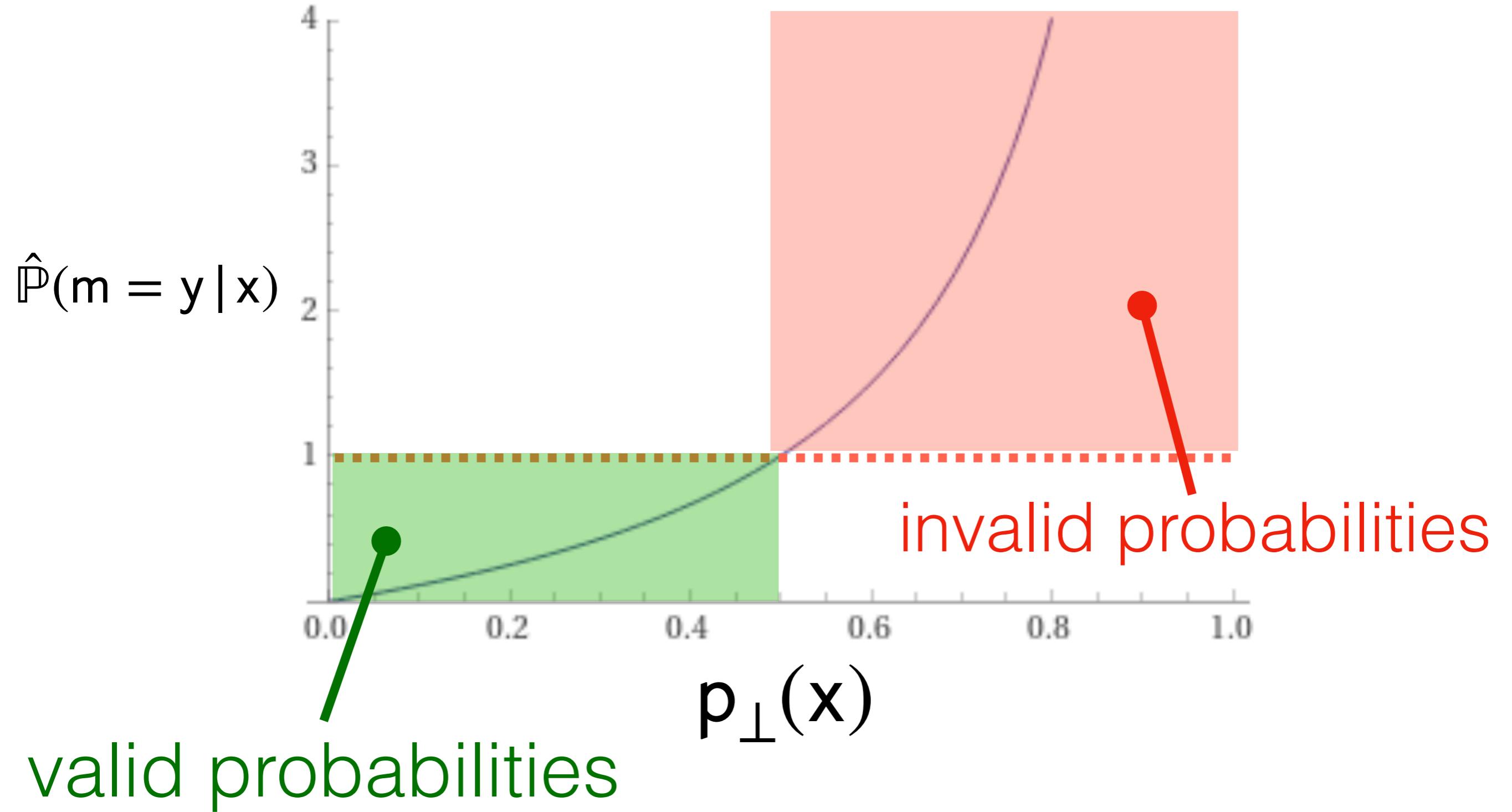
Estimating $P(m = y | x)$



Estimating $P(m = y | x)$



Estimating $P(m = y | x)$

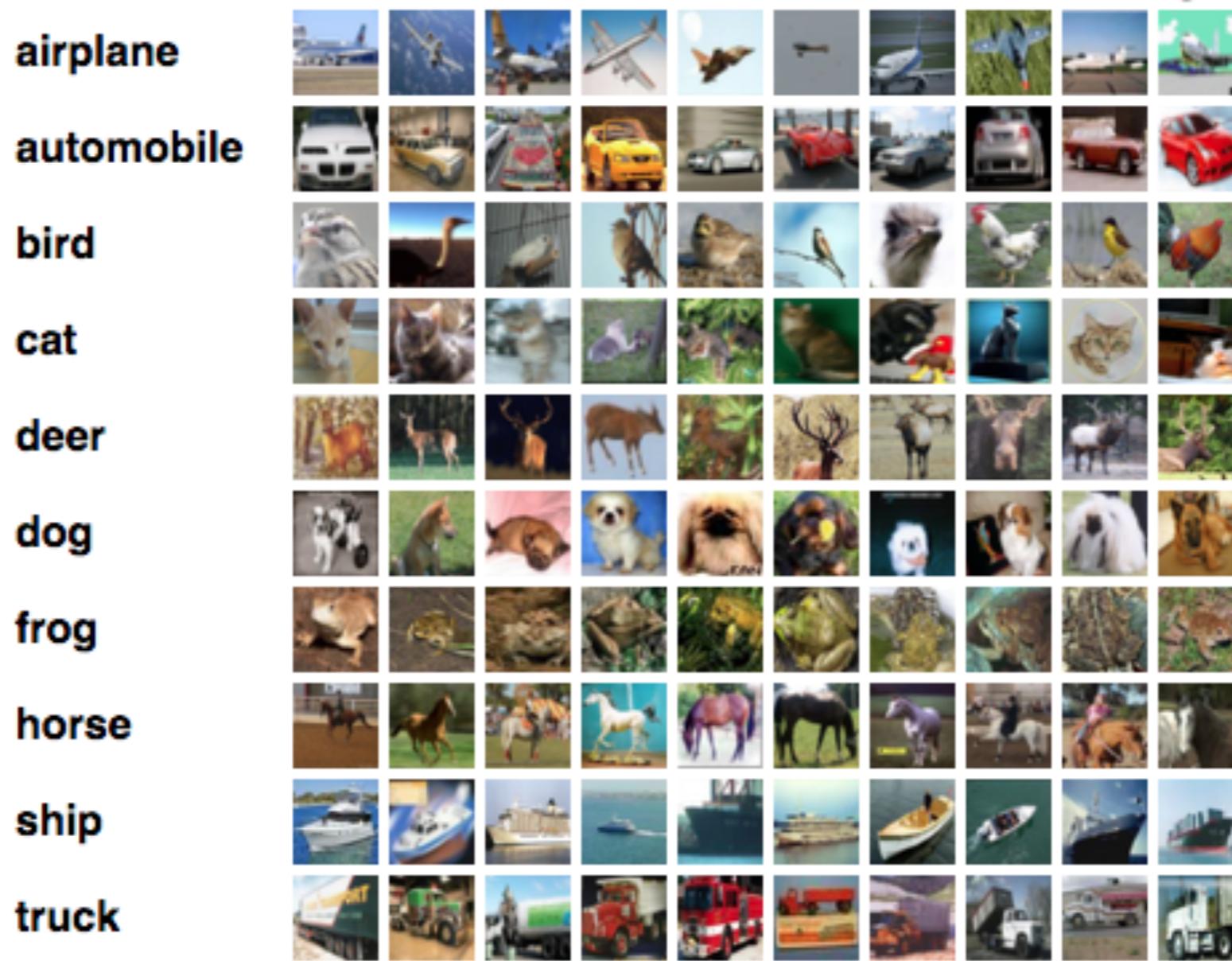


What happens in practice?

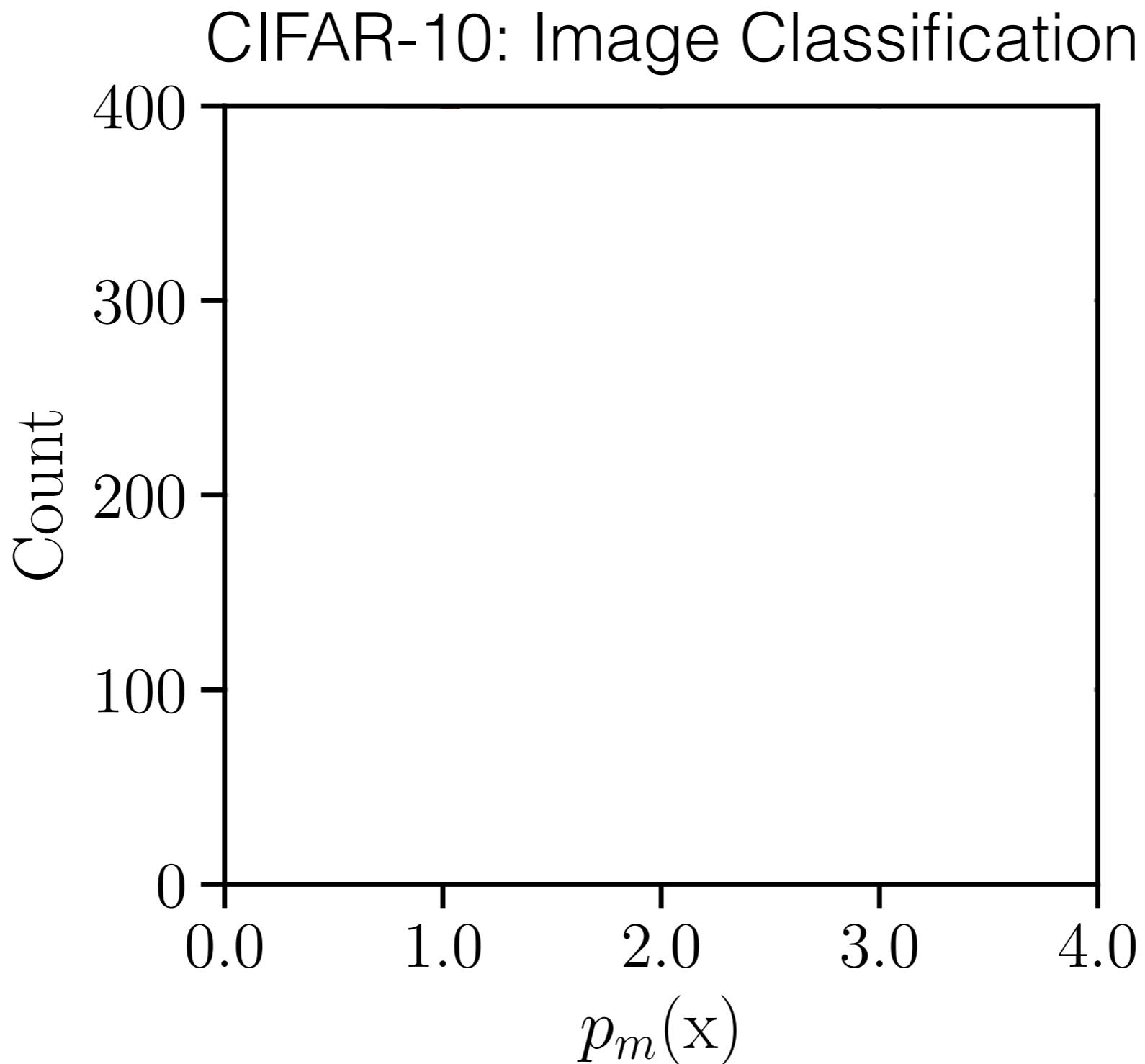
CIFAR-10: Image Classification

What happens in practice?

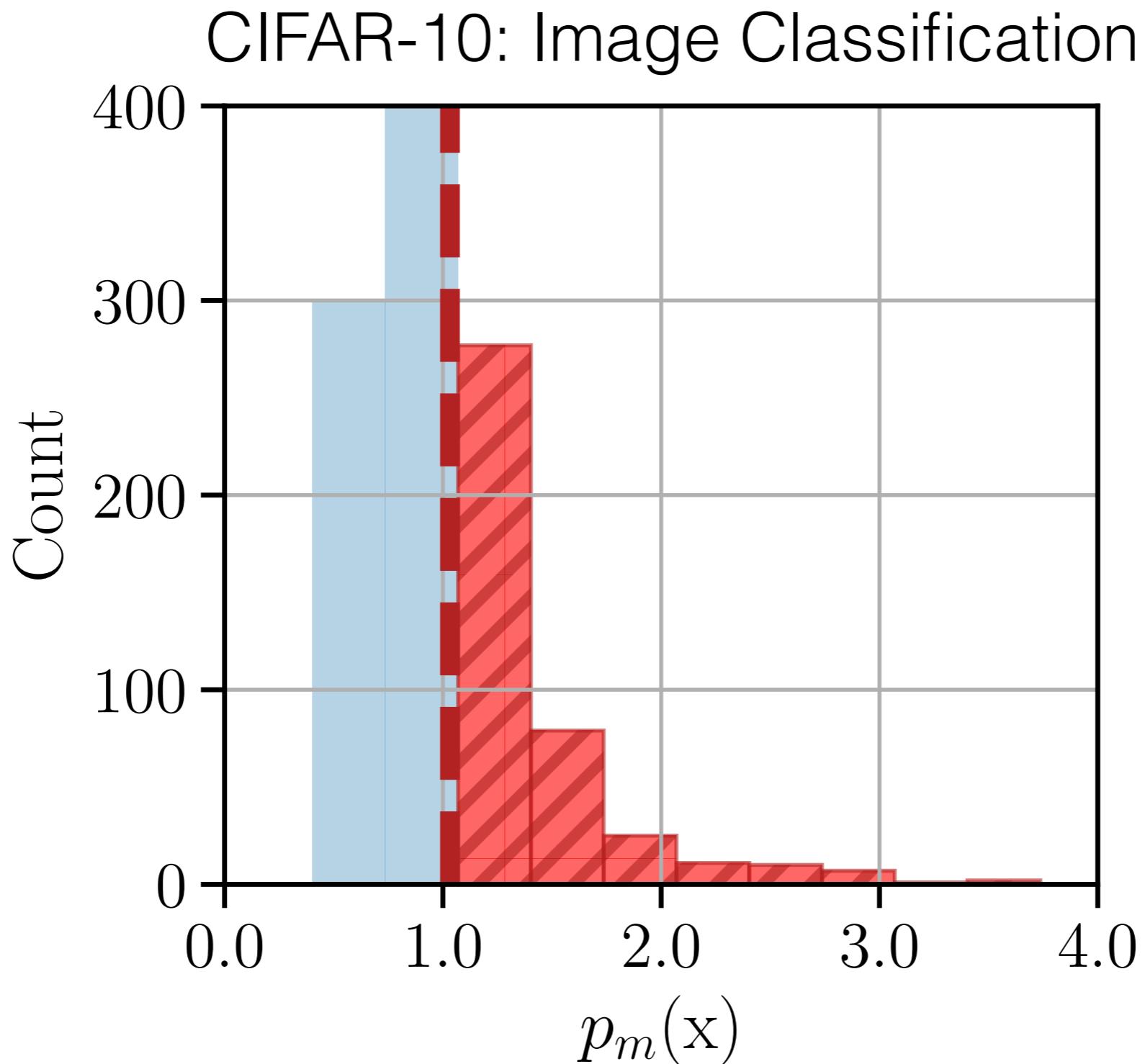
CIFAR-10: Image Classification



What happens in practice?



What happens in practice?

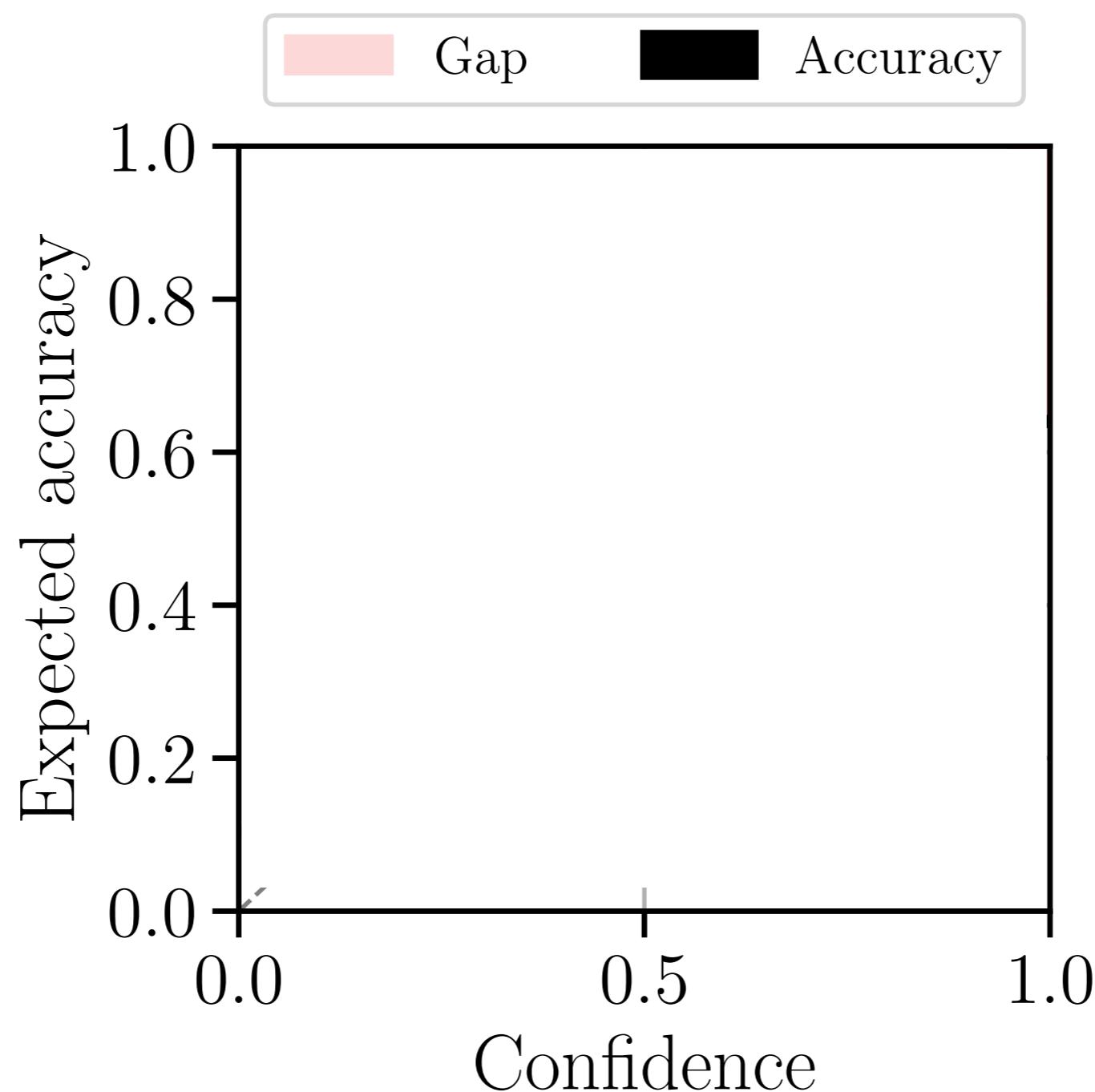


Truncation: $p_m(x) \in (0,1]$

CIFAR-10: Image Classification

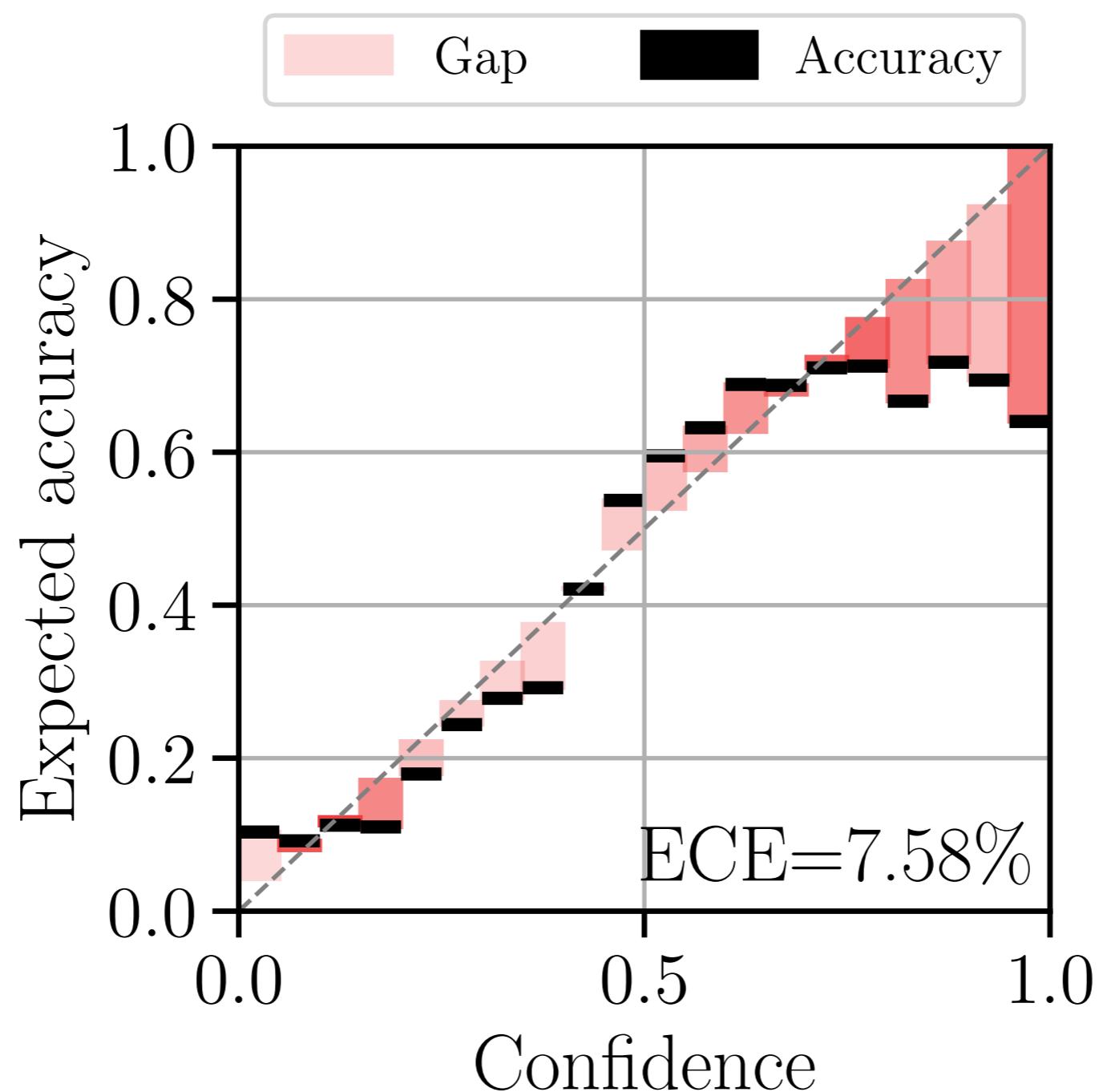
Truncation: $p_m(x) \in (0,1]$

CIFAR-10: Image Classification



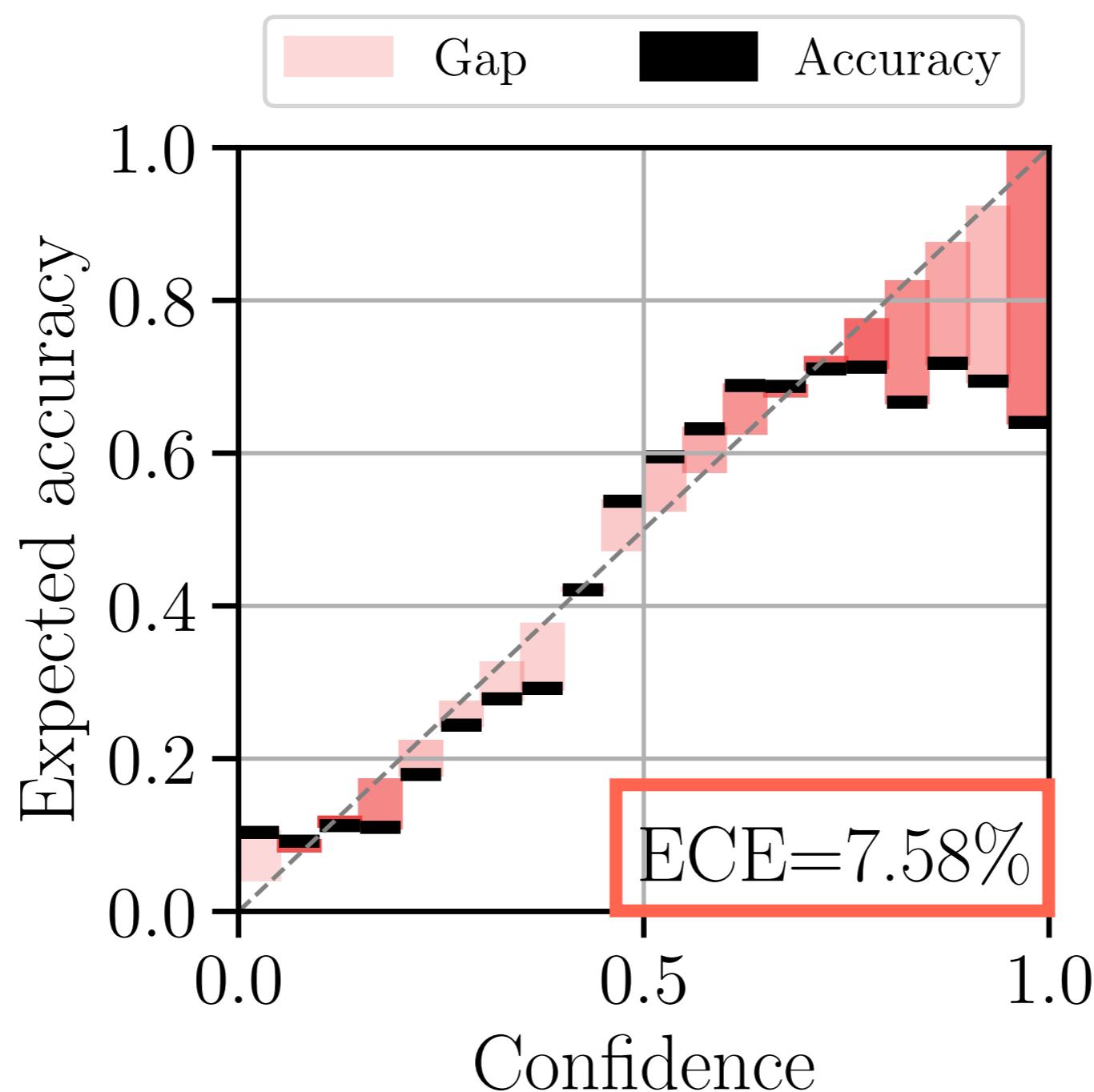
Truncation: $p_m(x) \in (0,1]$

CIFAR-10: Image Classification



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

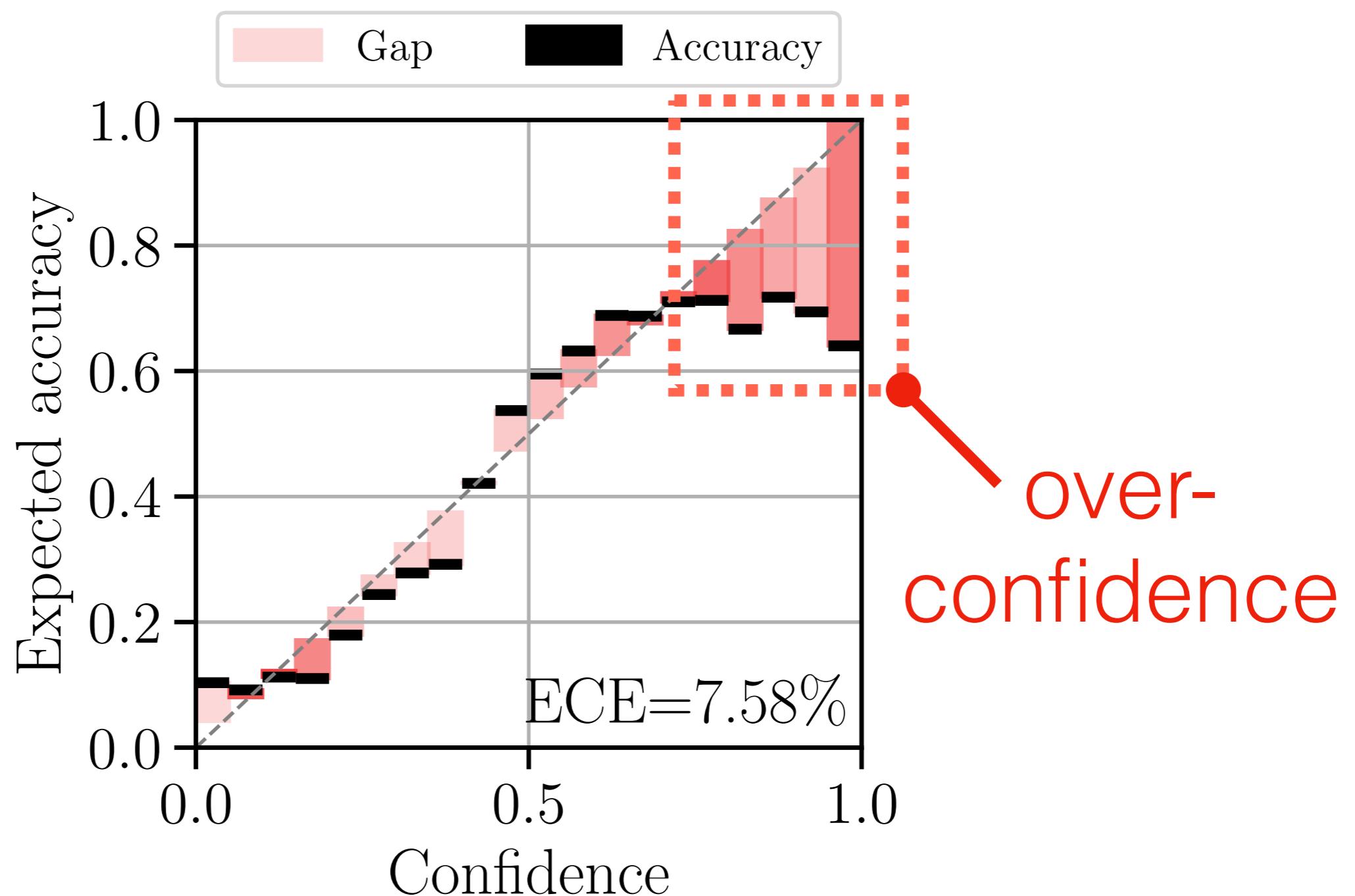
CIFAR-10: Image Classification



$ECE =$
expected
calibration
error

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

CIFAR-10: Image Classification

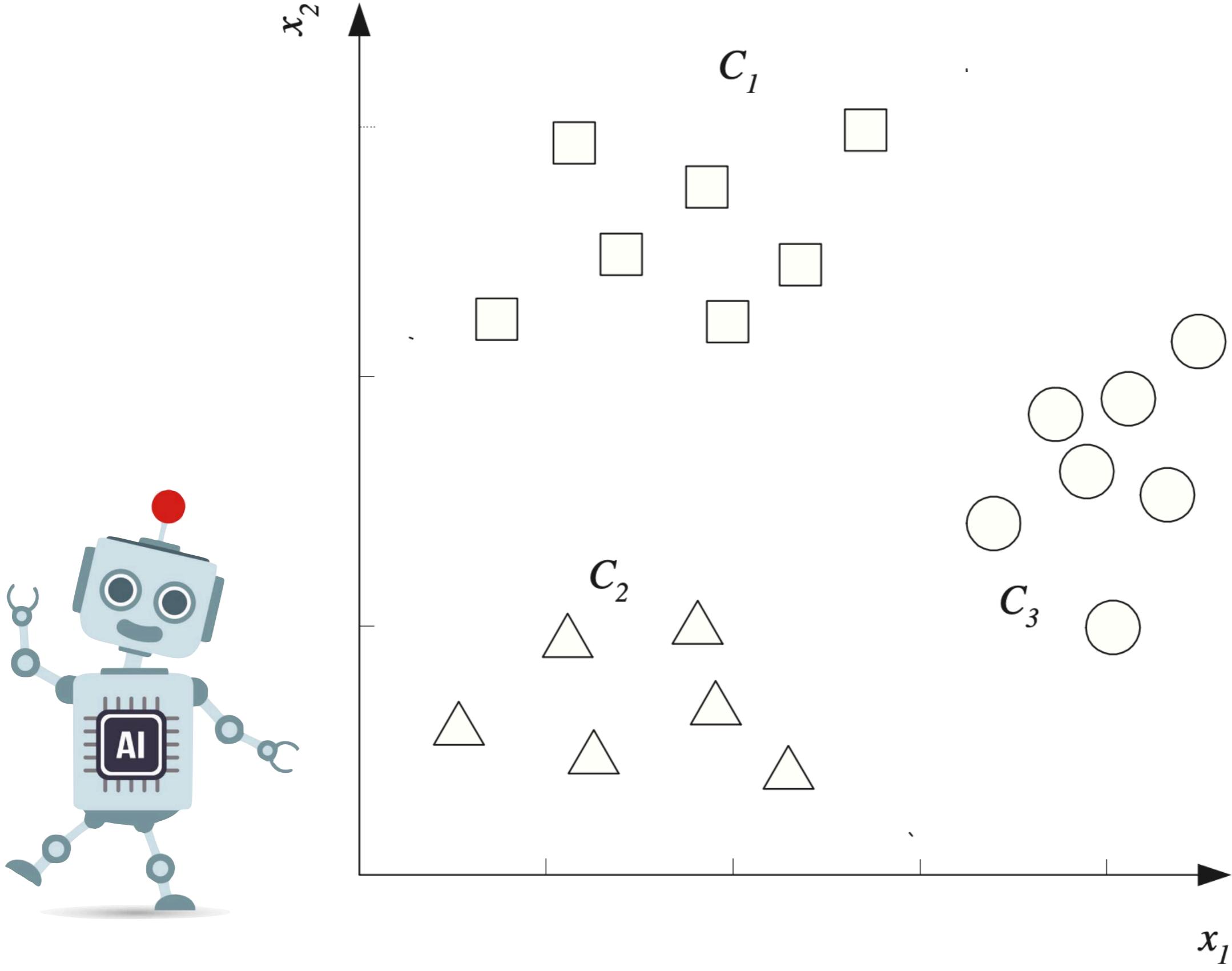


Solution

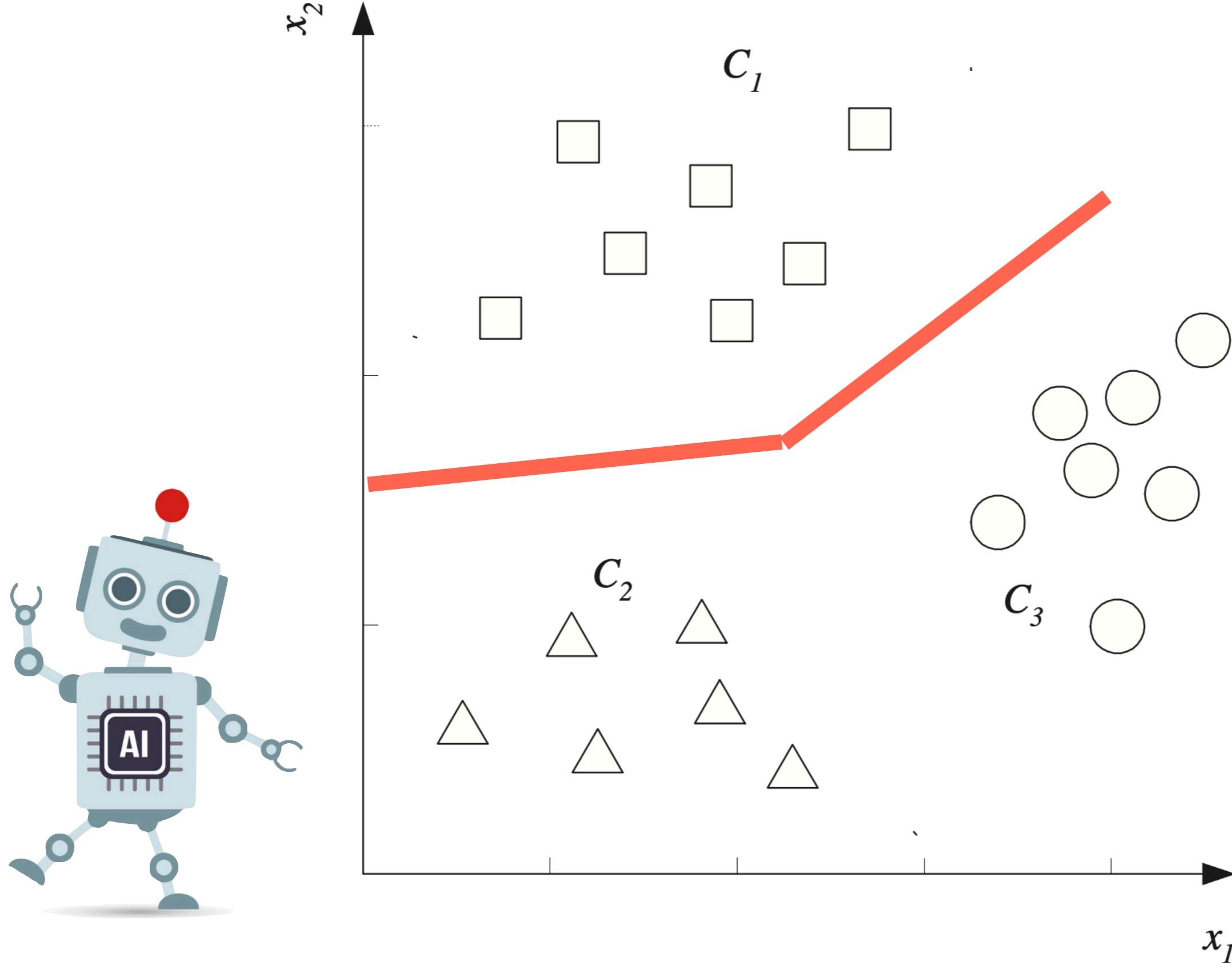
Learning to Defer
using a One-vs-All Classifier

Solution

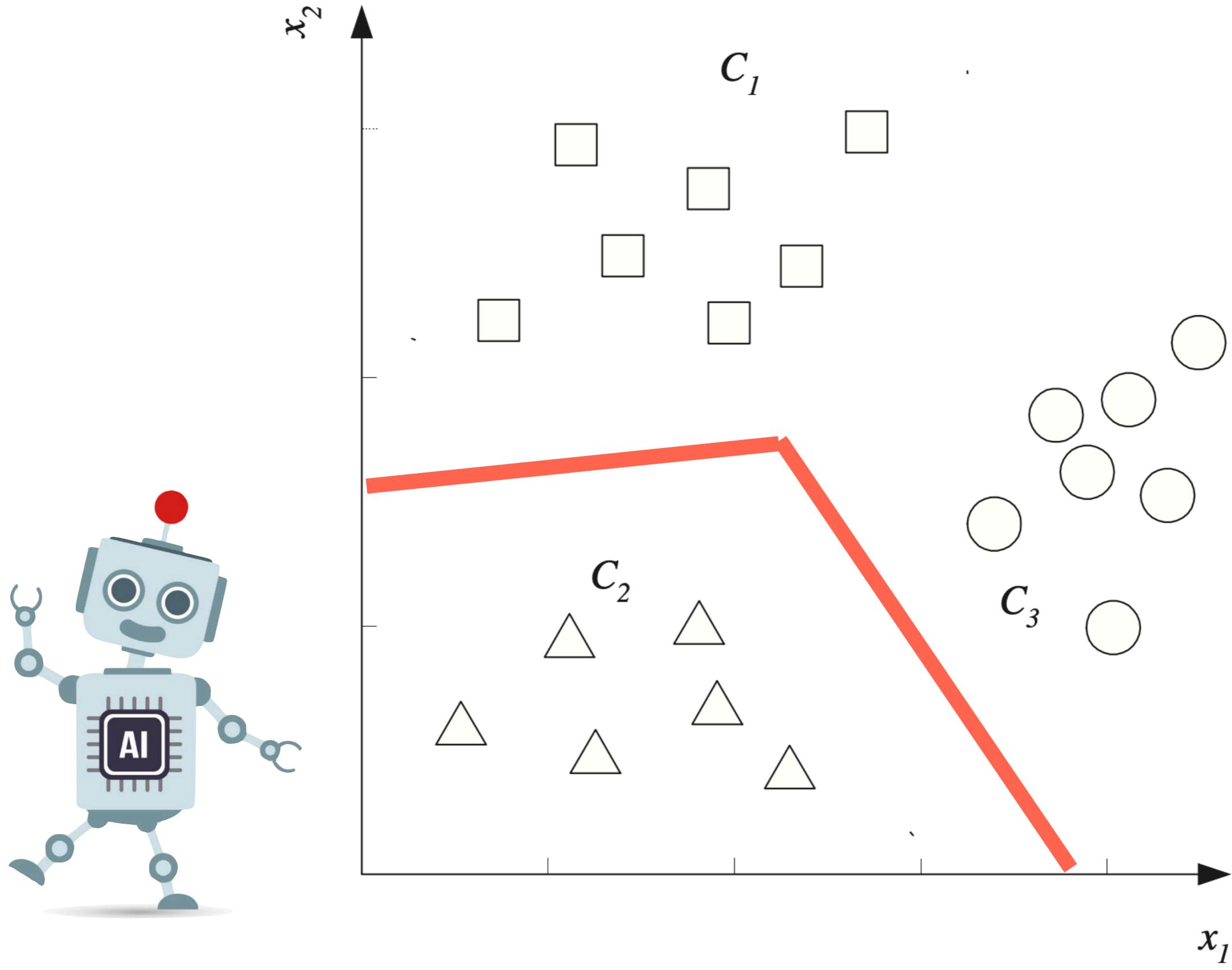
Learning to Defer
using a One-vs-All Classifier



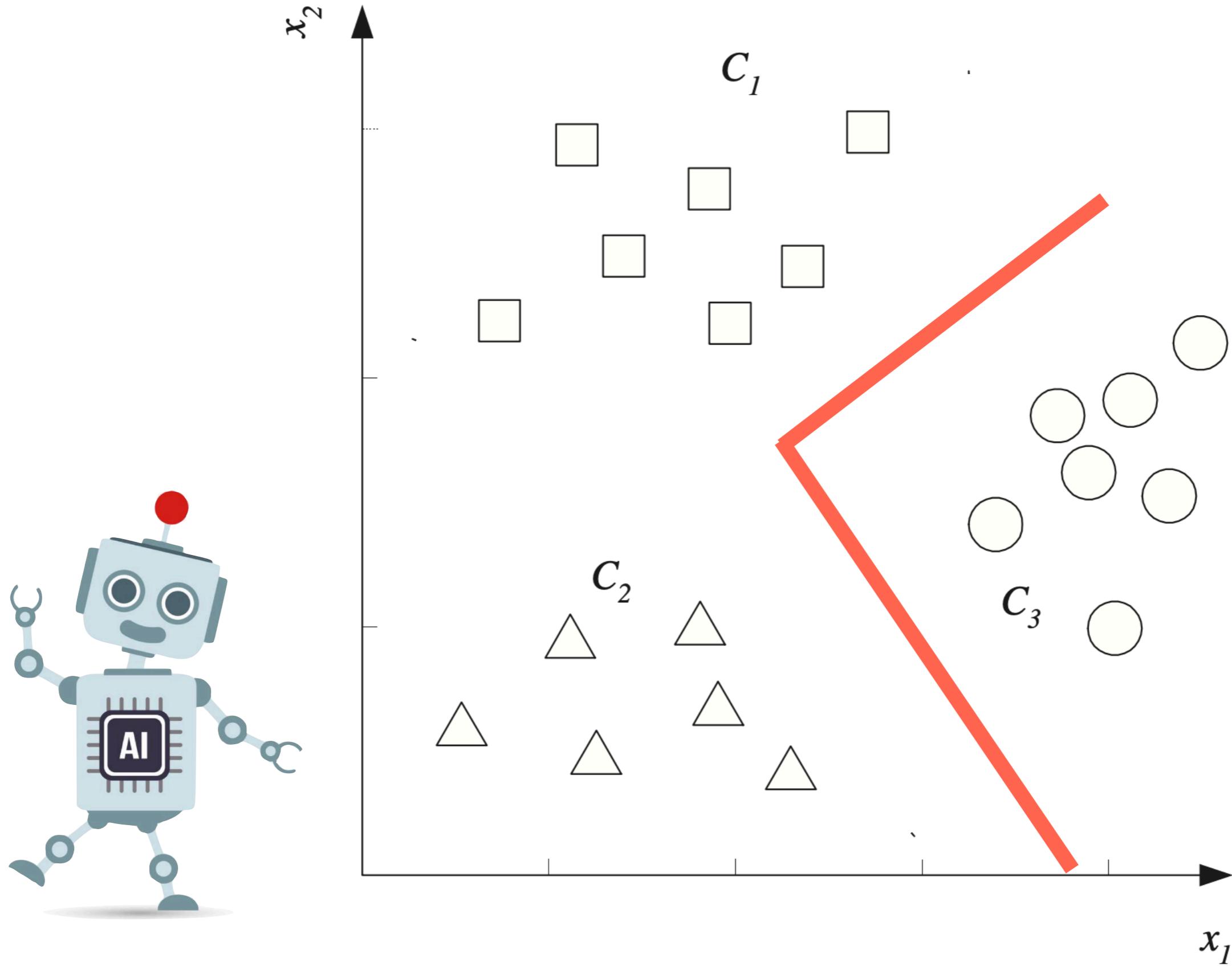
One-vs-All Classifier



One-vs-All Classifier

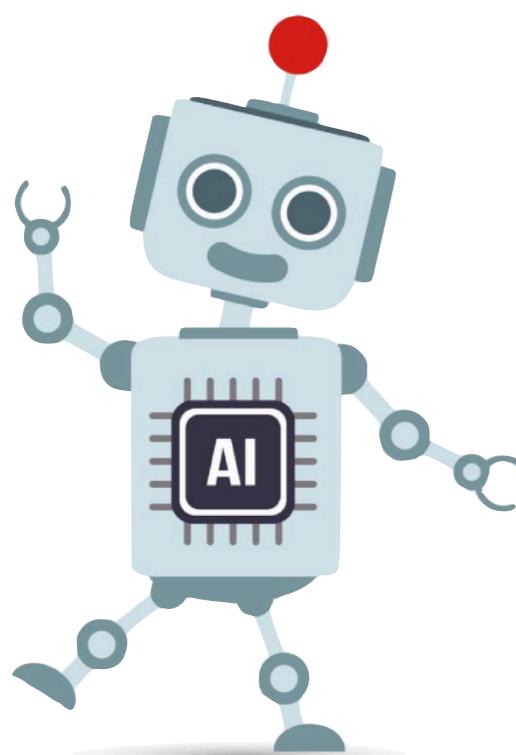


One-vs-All Classifier

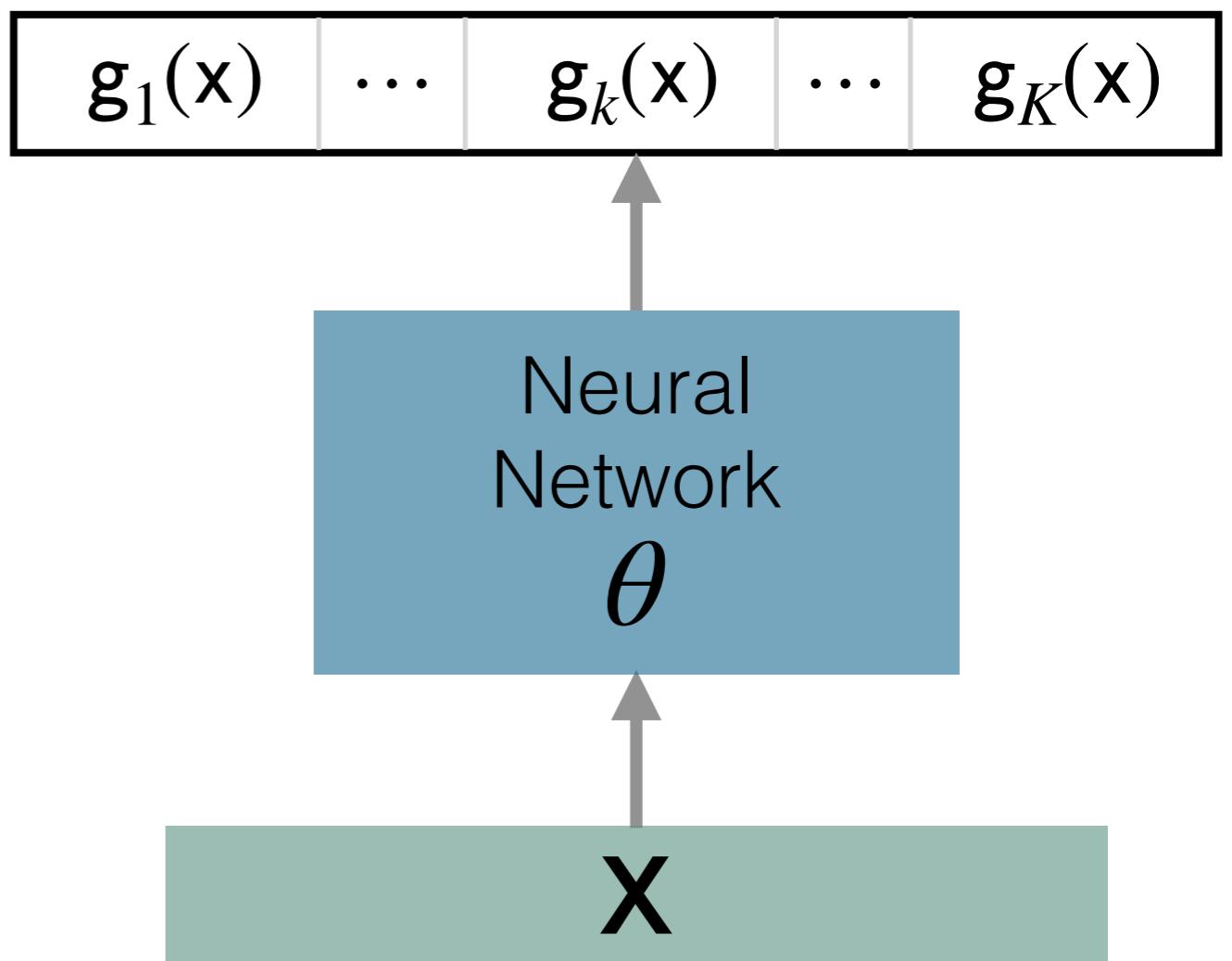


One-vs-All Classifier

Goal: For input features x ,
predict membership in
one of K classes.

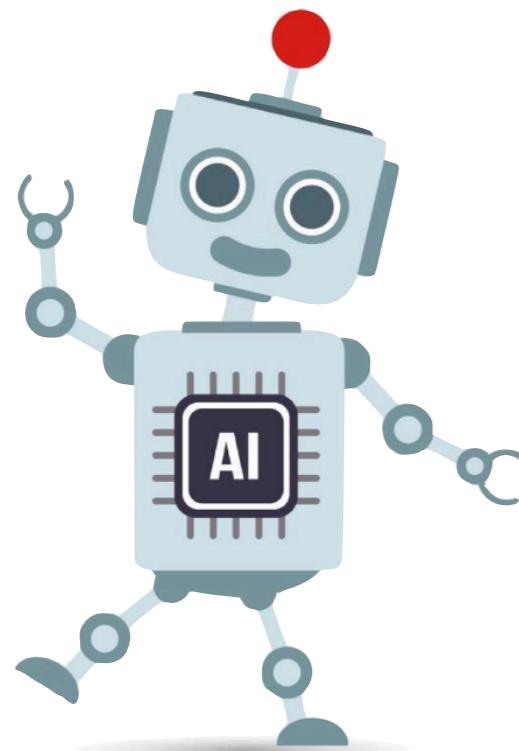


One-vs-All Classifier

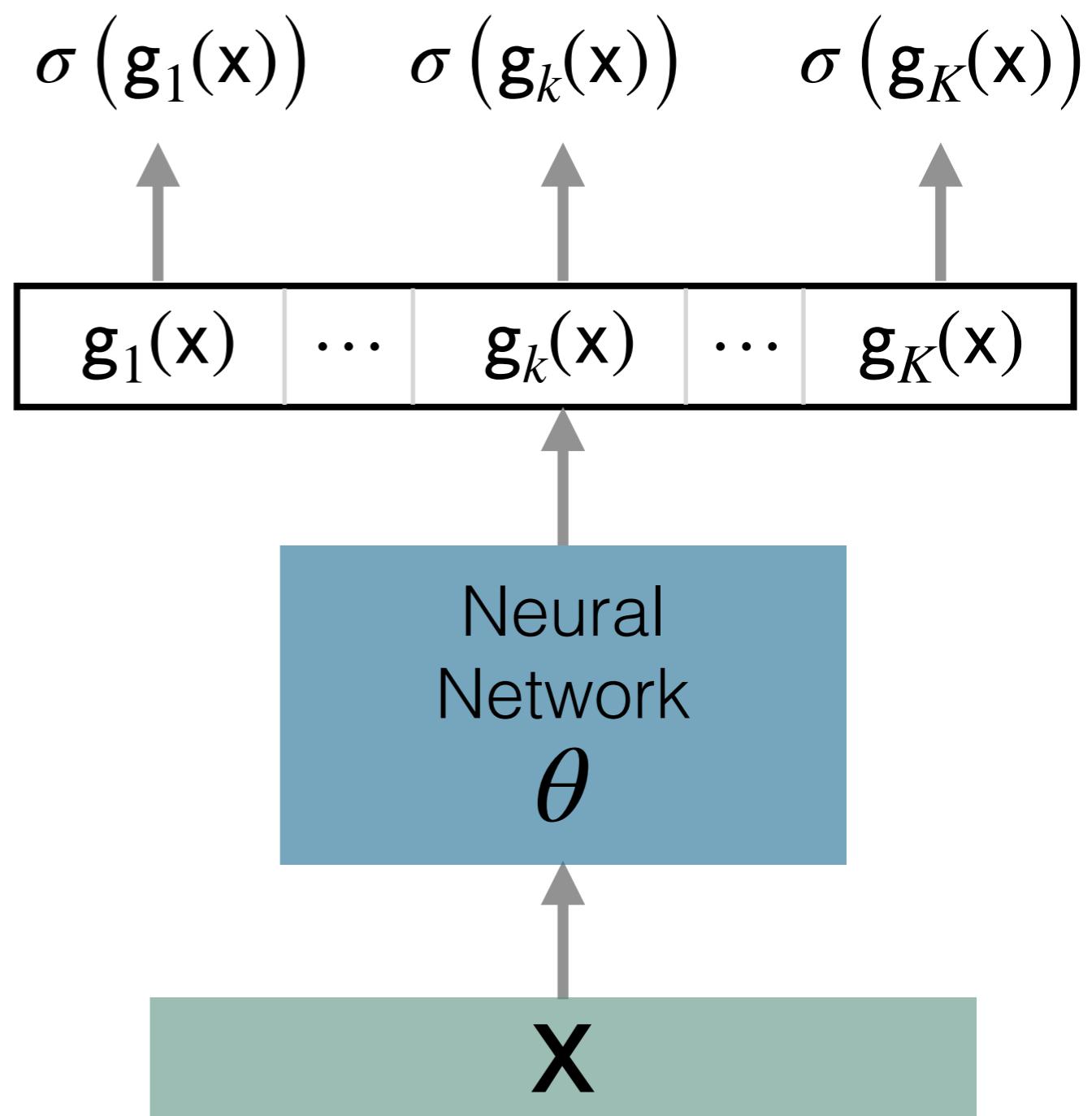


Goal: For input features \mathbf{x} ,
predict membership in
one of K classes.

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

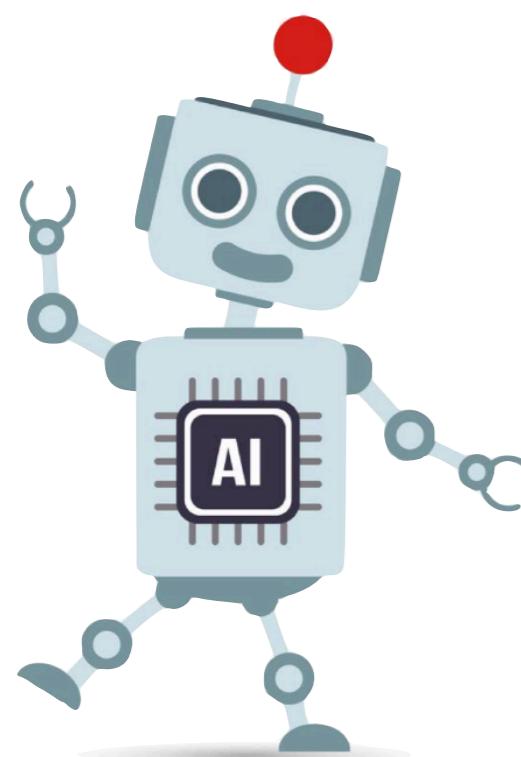


One-vs-All Classifier

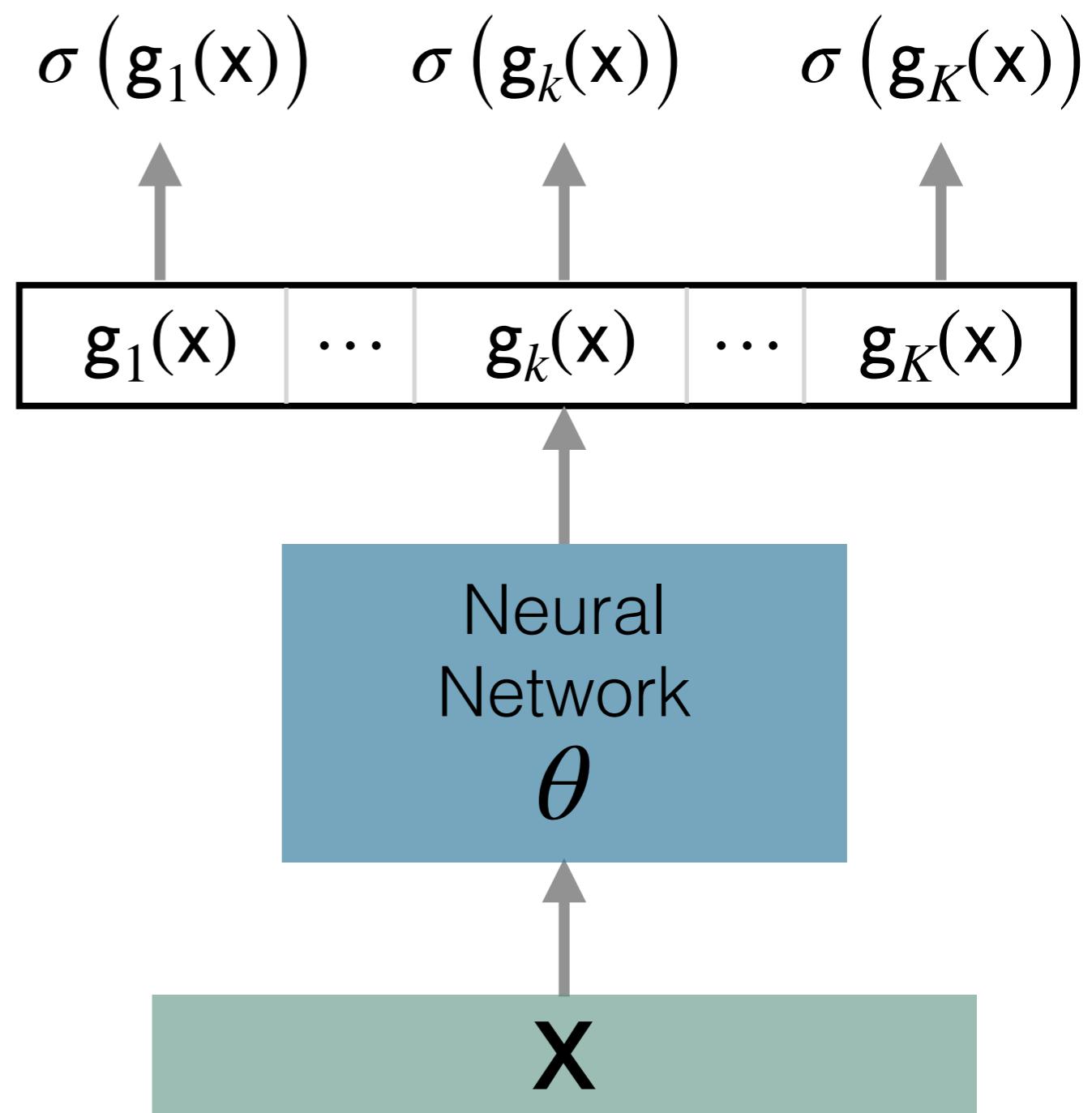


Goal: For input features \mathbf{x} ,
predict membership in
one of K classes.

$$P(y | \mathbf{x}) = \sigma(g_y(\mathbf{x}))$$

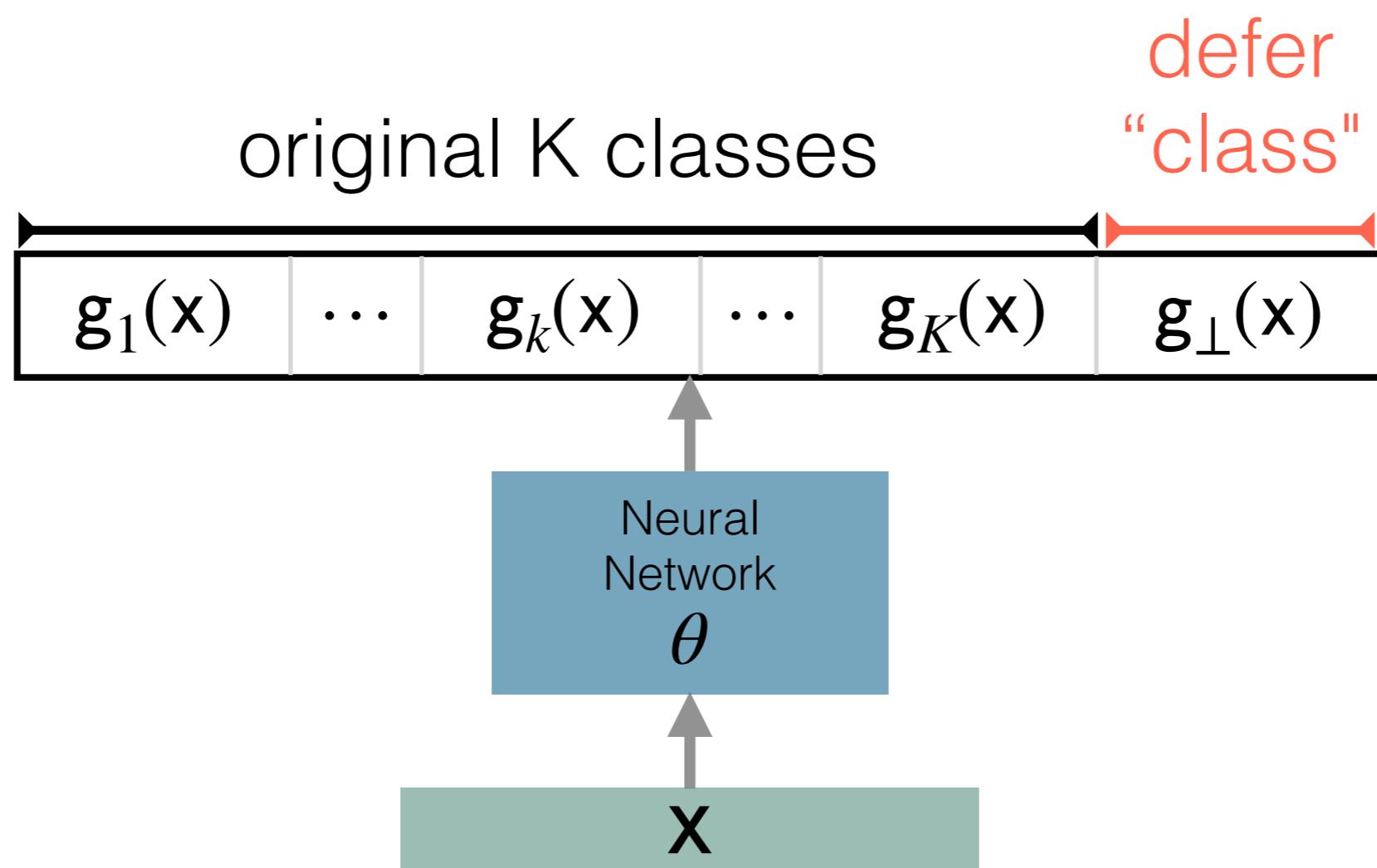


One-vs-All Classifier

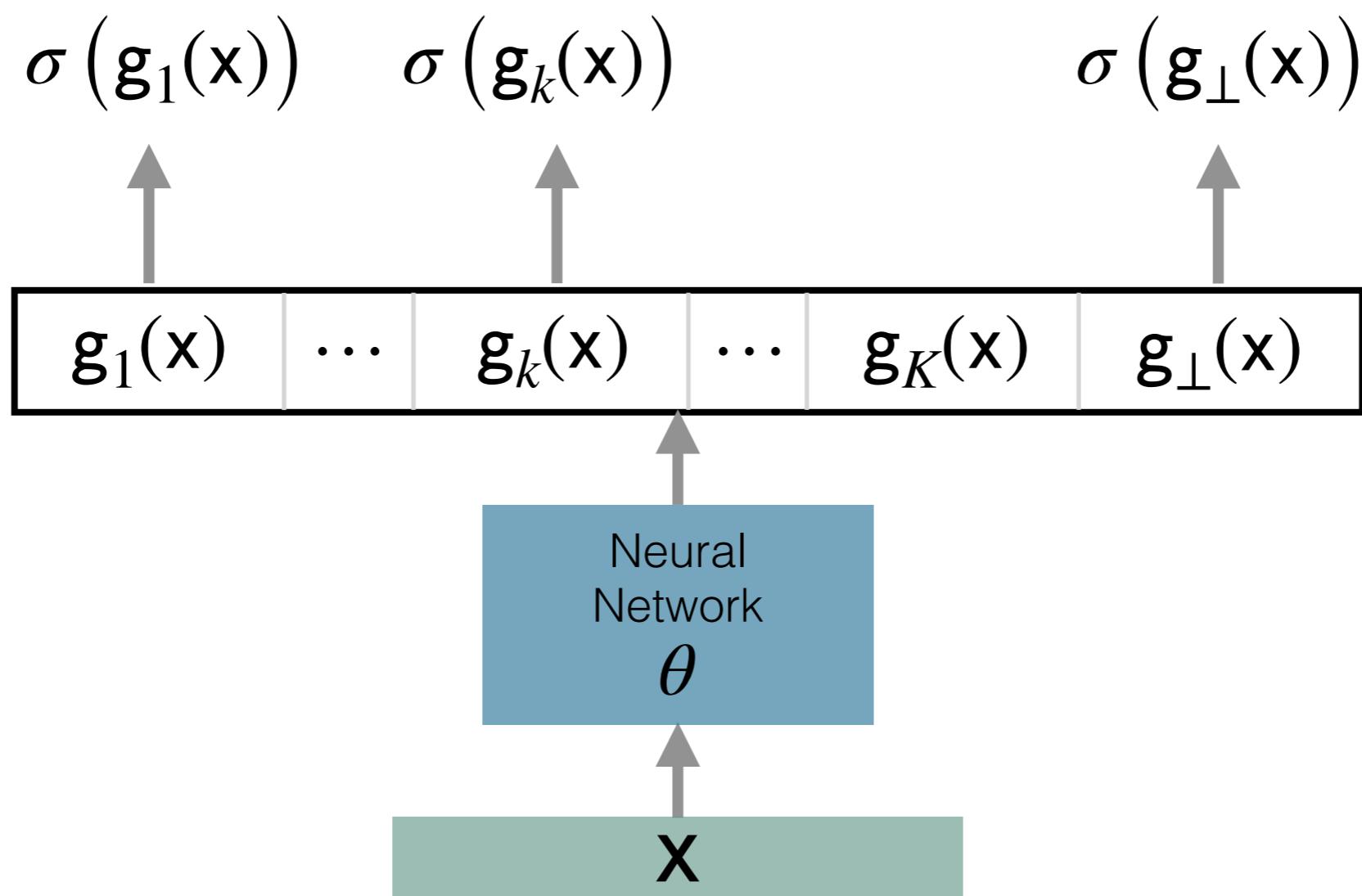


Learning-to-Defer: One-vs-All Parameterization

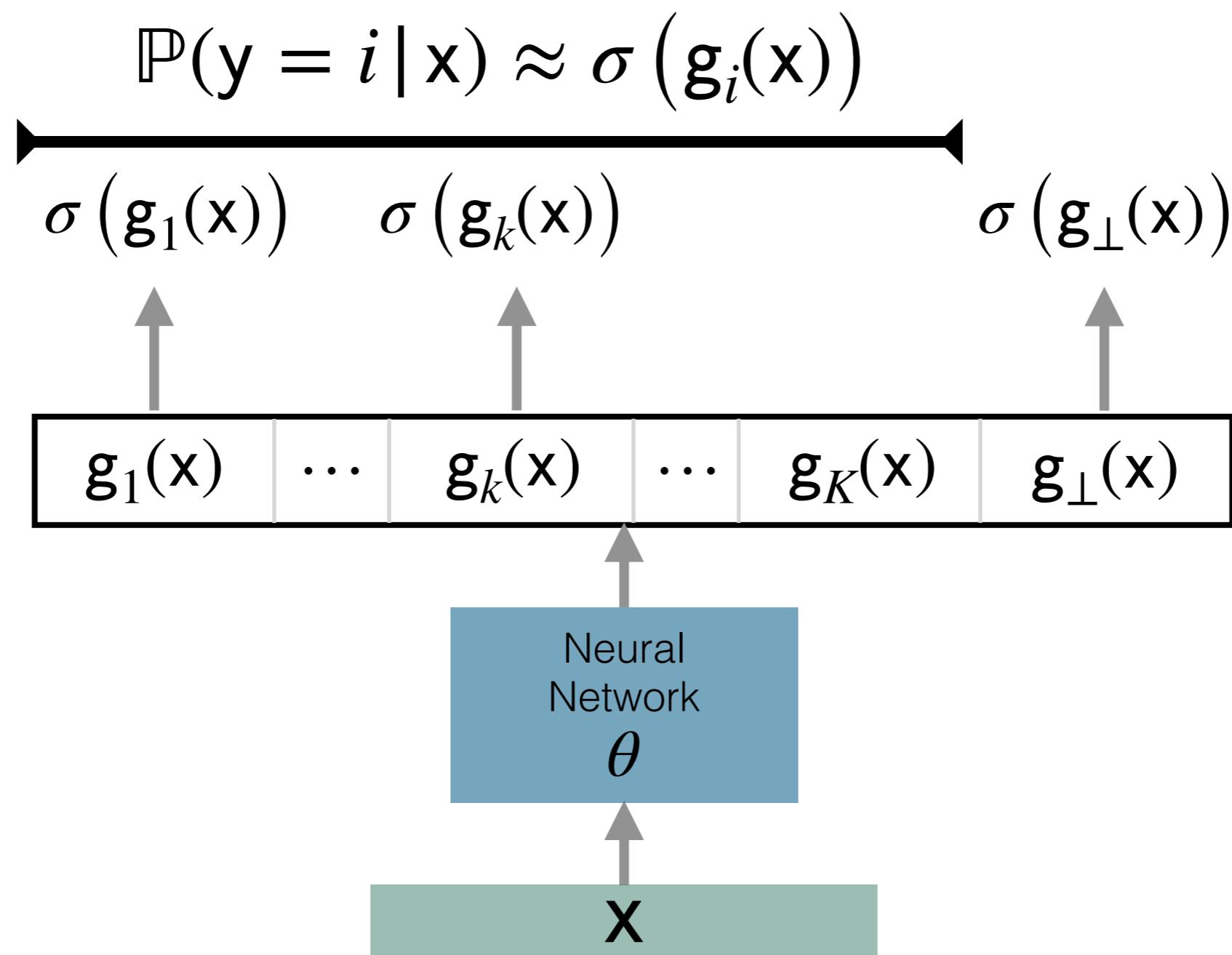
Learning-to-Defer: One-vs-All Parameterization



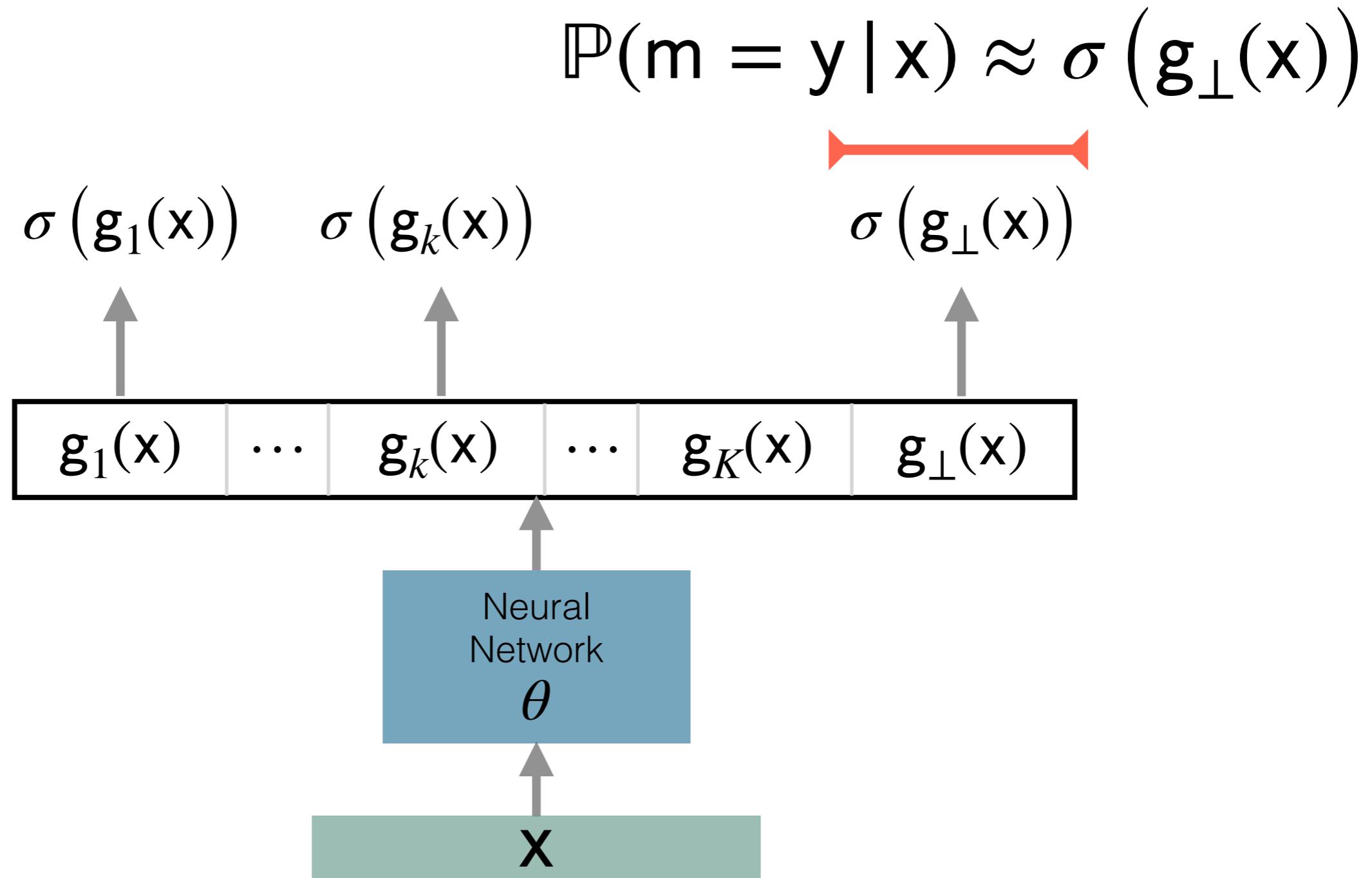
Learning-to-Defer: One-vs-All Parameterization



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



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



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

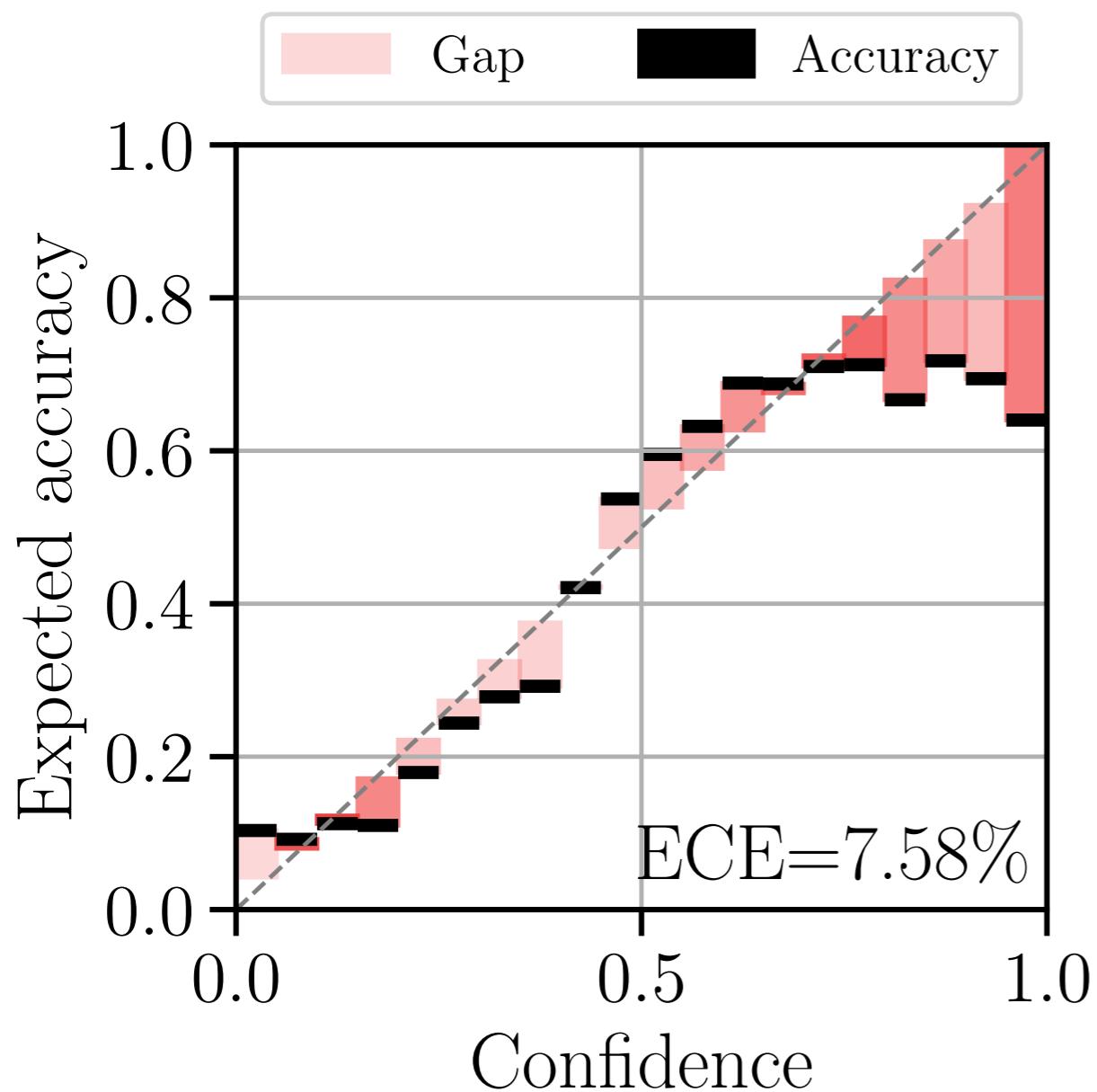
Theorem 4.1 + Corollary 4.2 [Verma & Nalisnick, 2022]:
The one-vs-all loss is a *consistent* surrogate for
the 0-1 learning-to-defer loss.

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

CIFAR-10: Softmax

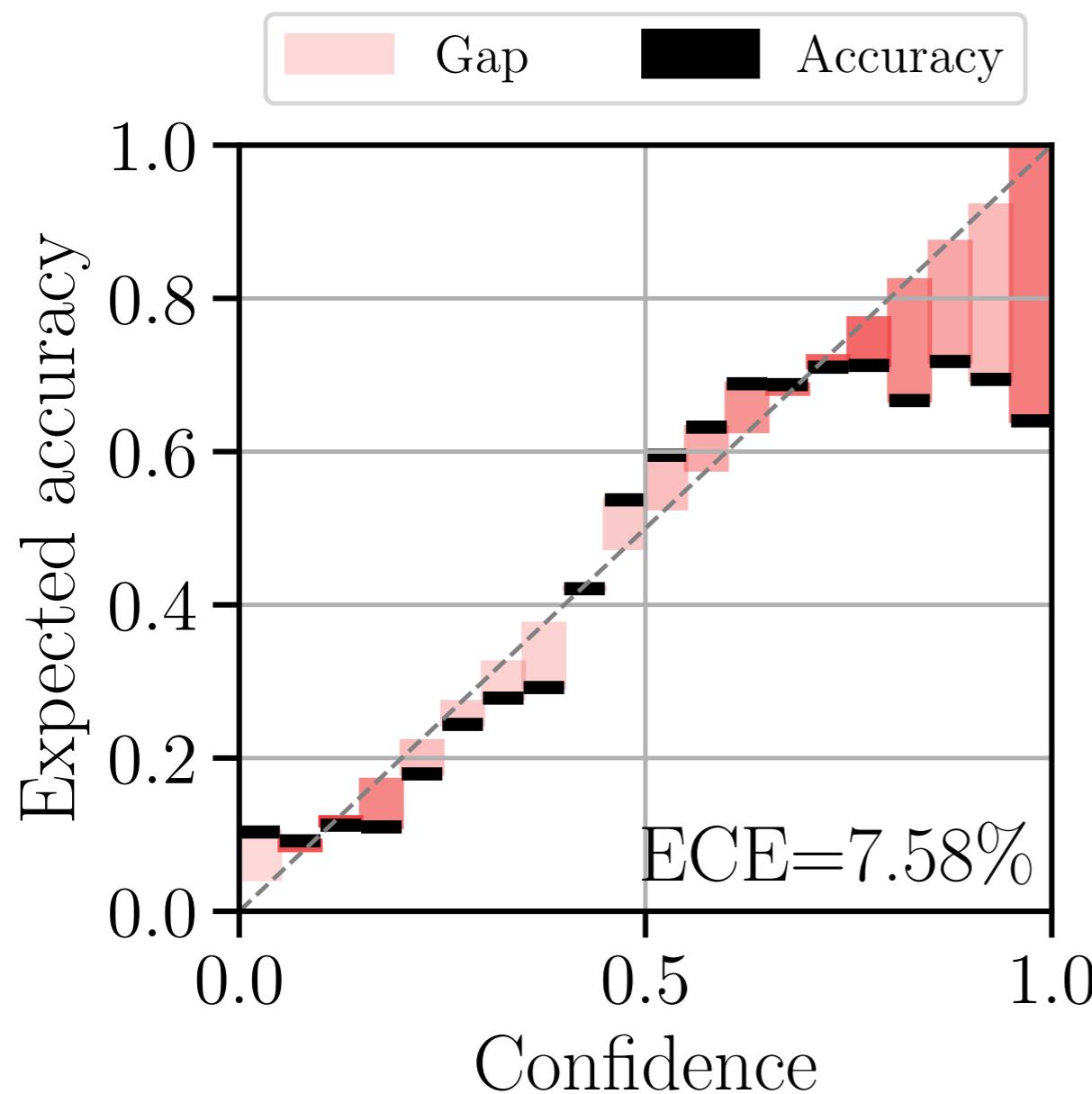
CIFAR-10: One-vs-All

CIFAR-10: Softmax

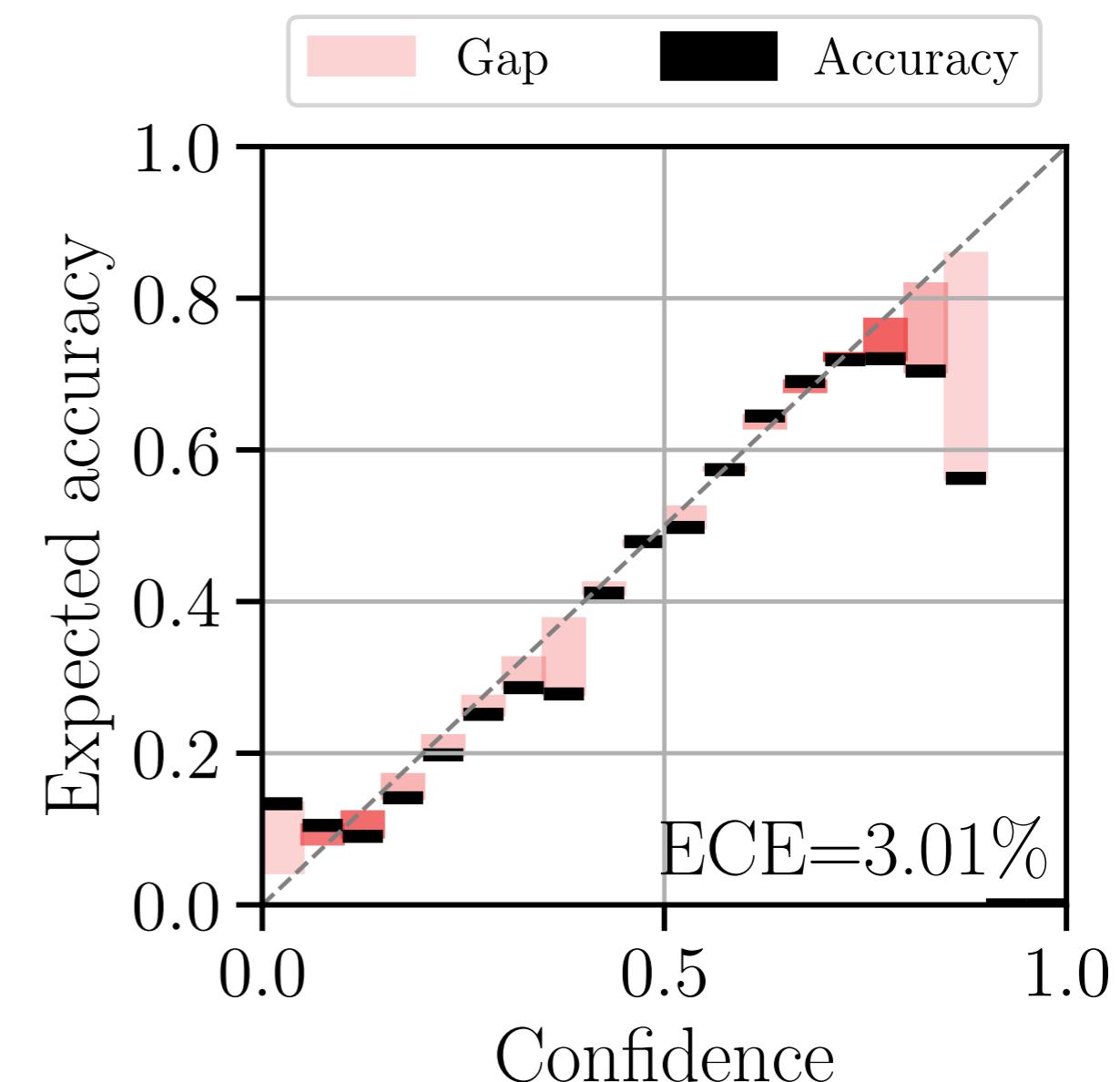


CIFAR-10: One-vs-All

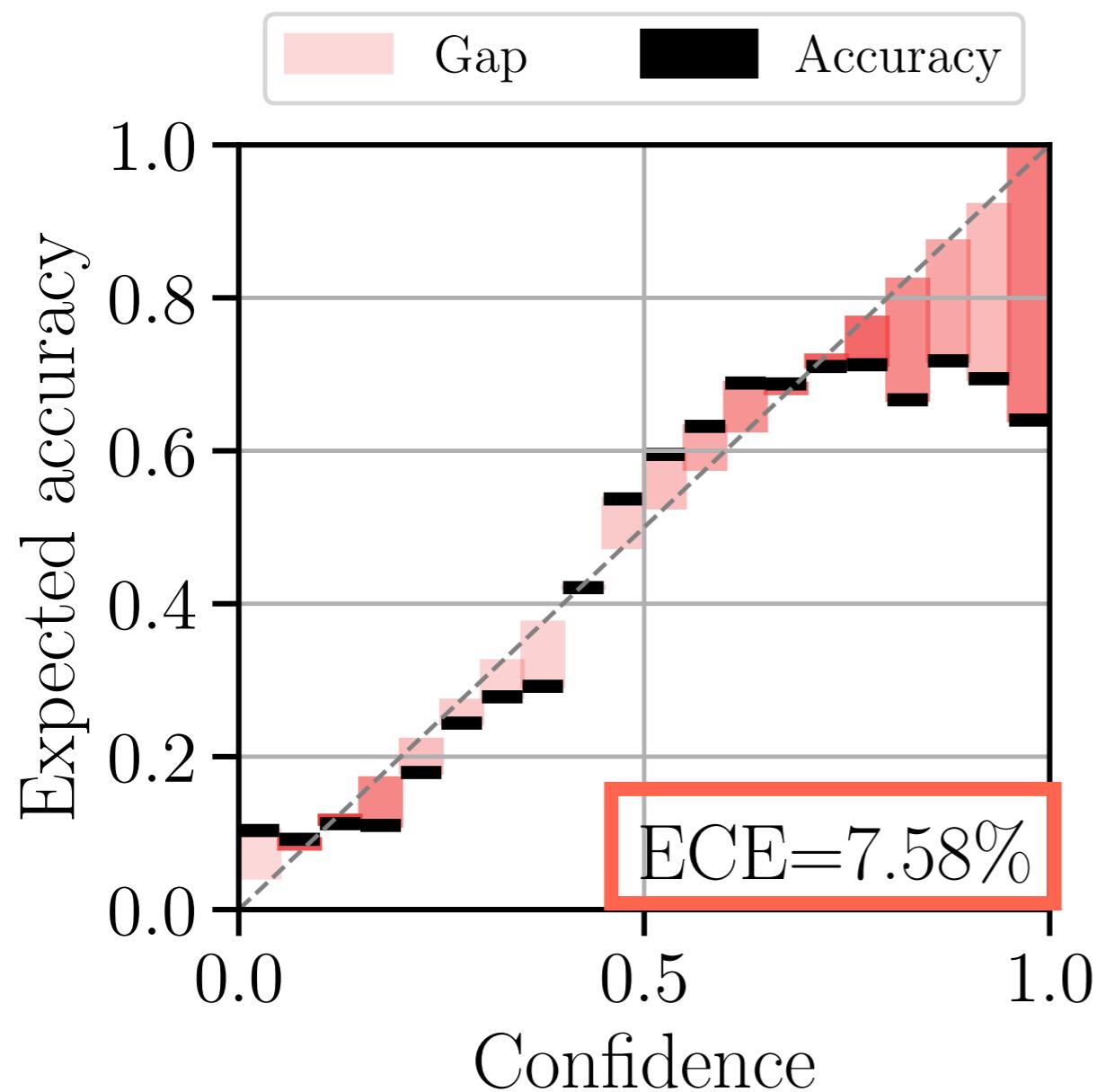
CIFAR-10: Softmax



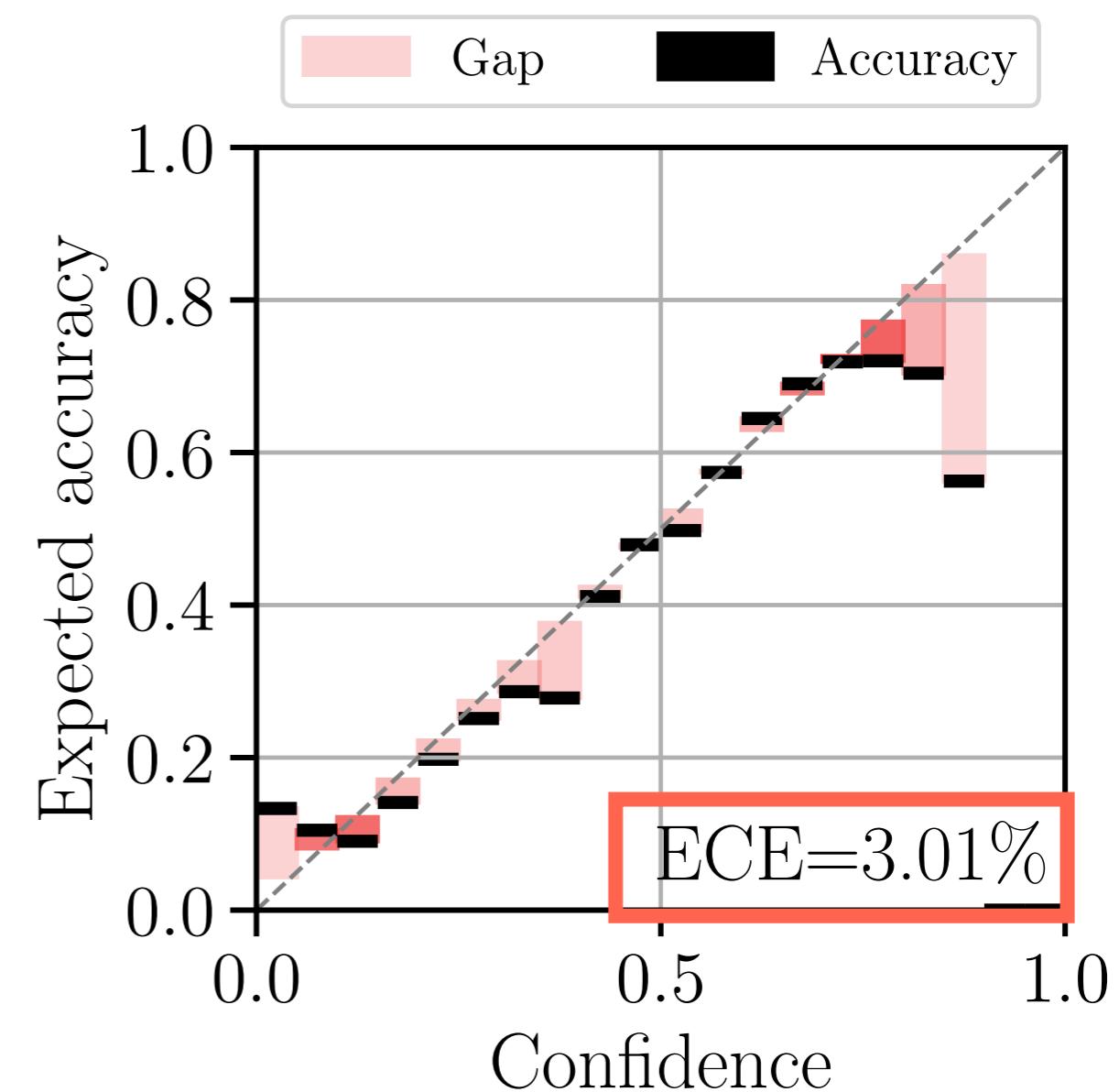
CIFAR-10: One-vs-All



CIFAR-10: Softmax

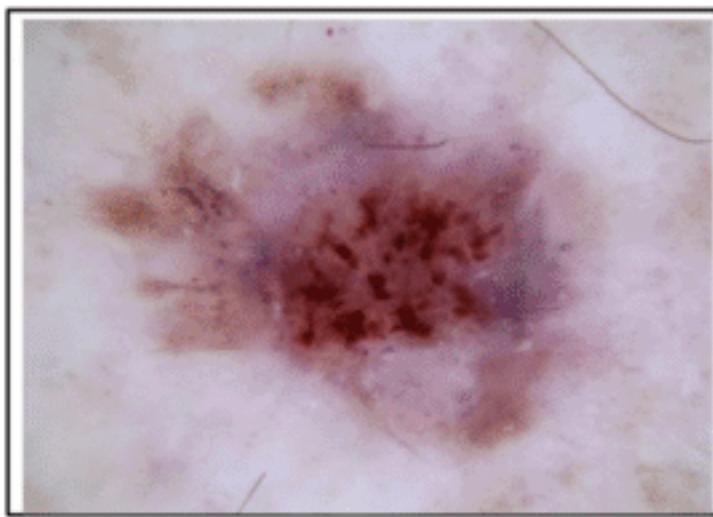


CIFAR-10: One-vs-All

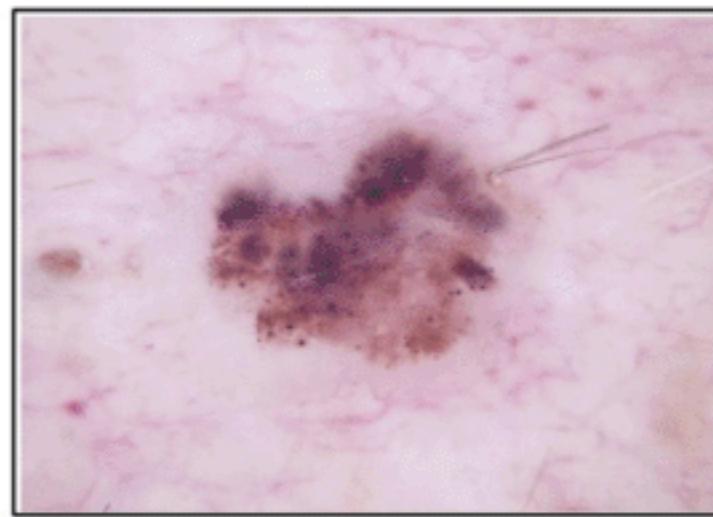


HAM10000: Skin Lesion Classification

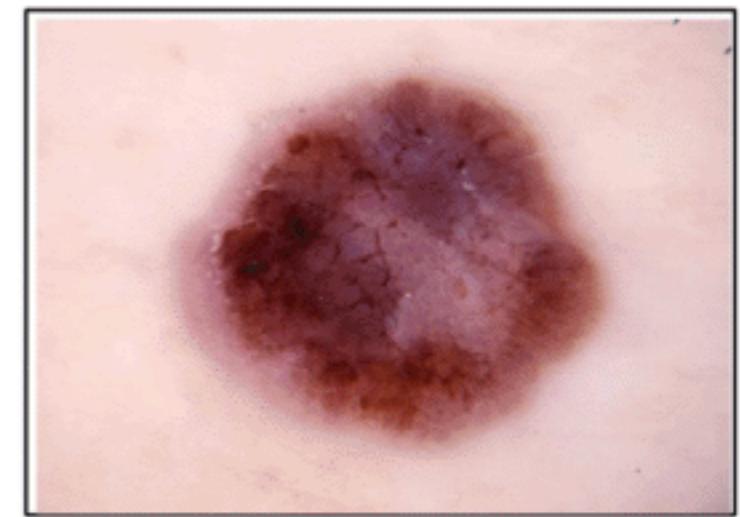
(AKIEC)



(BCC)



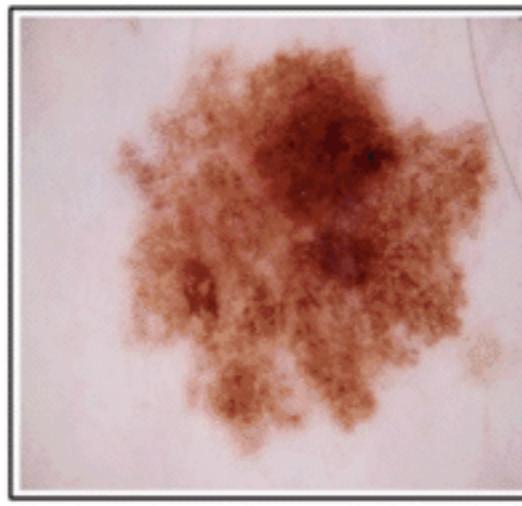
(BKL)



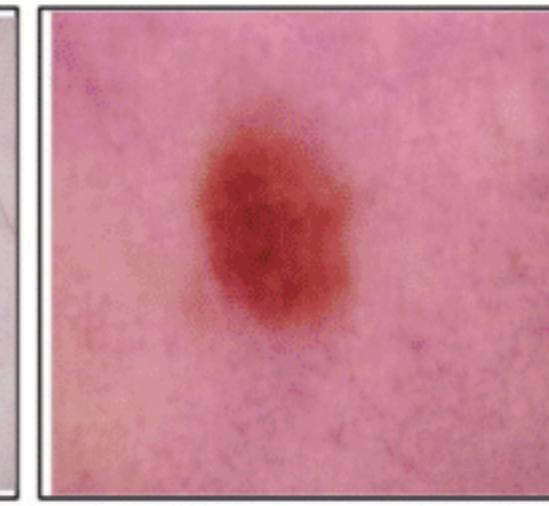
(DF)



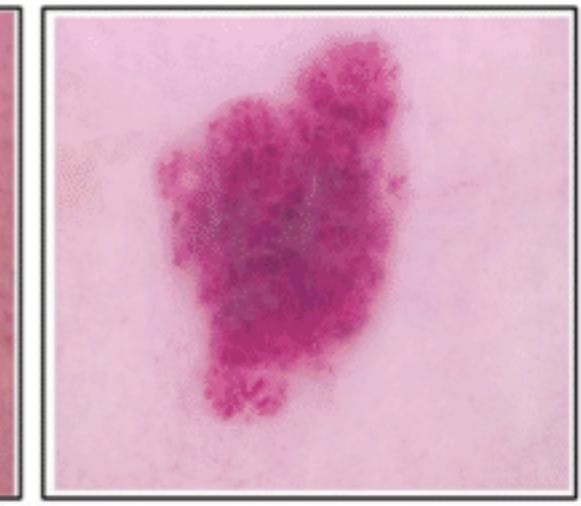
(MEL)



(NV)



(VASC)

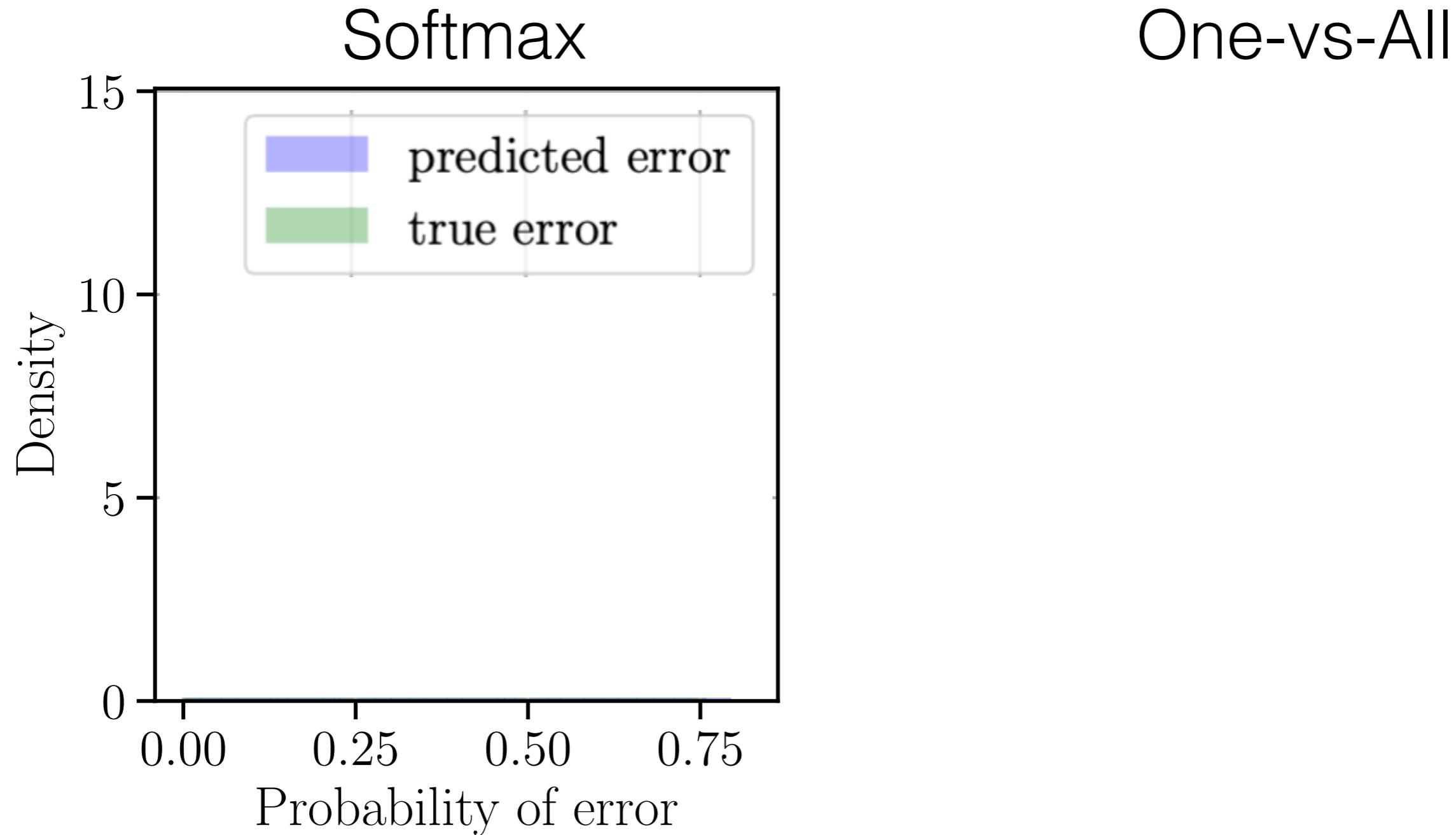


HAM10000: Skin Lesion Classification

Softmax

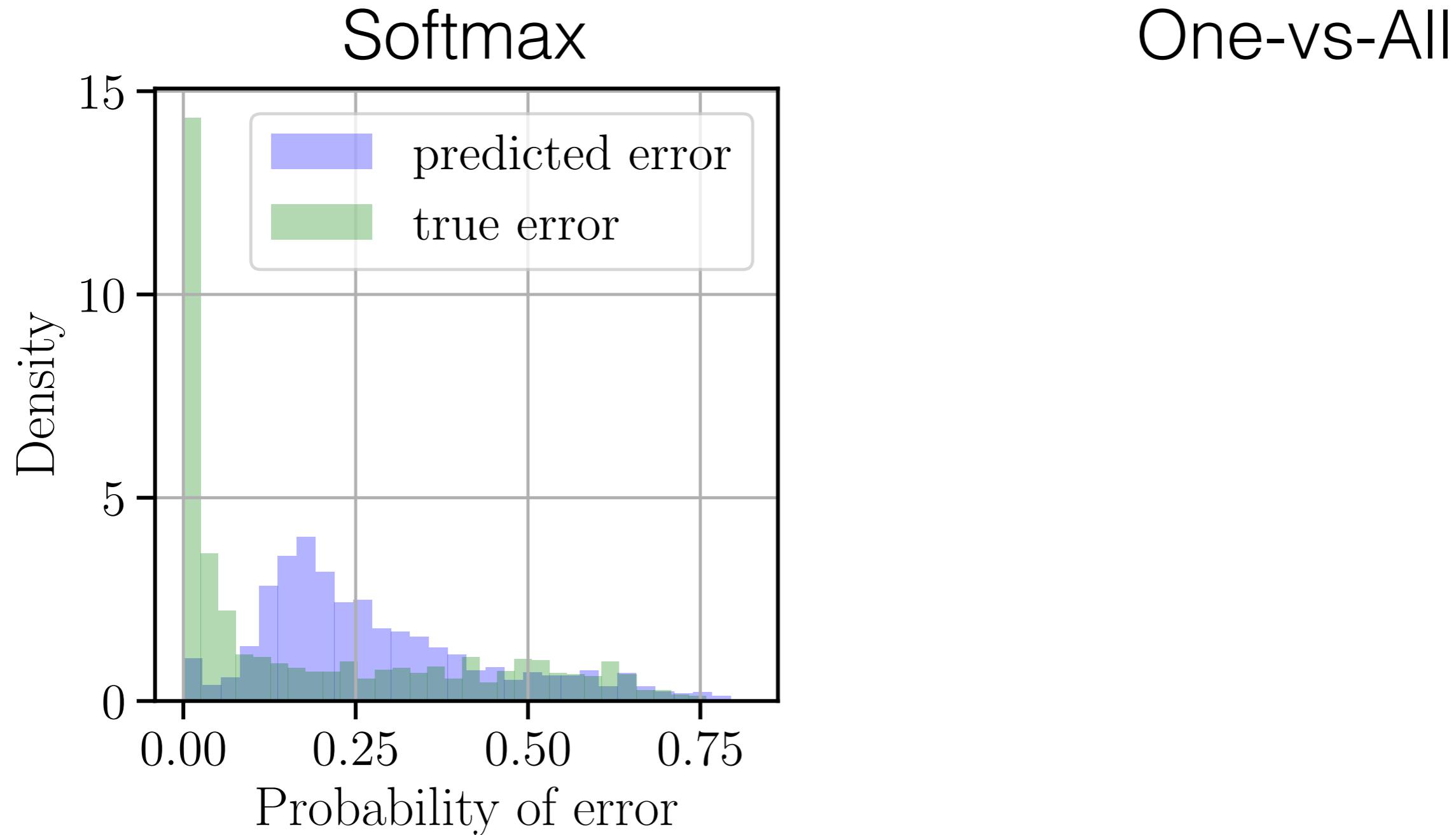
One-vs-All

HAM10000: Skin Lesion Classification



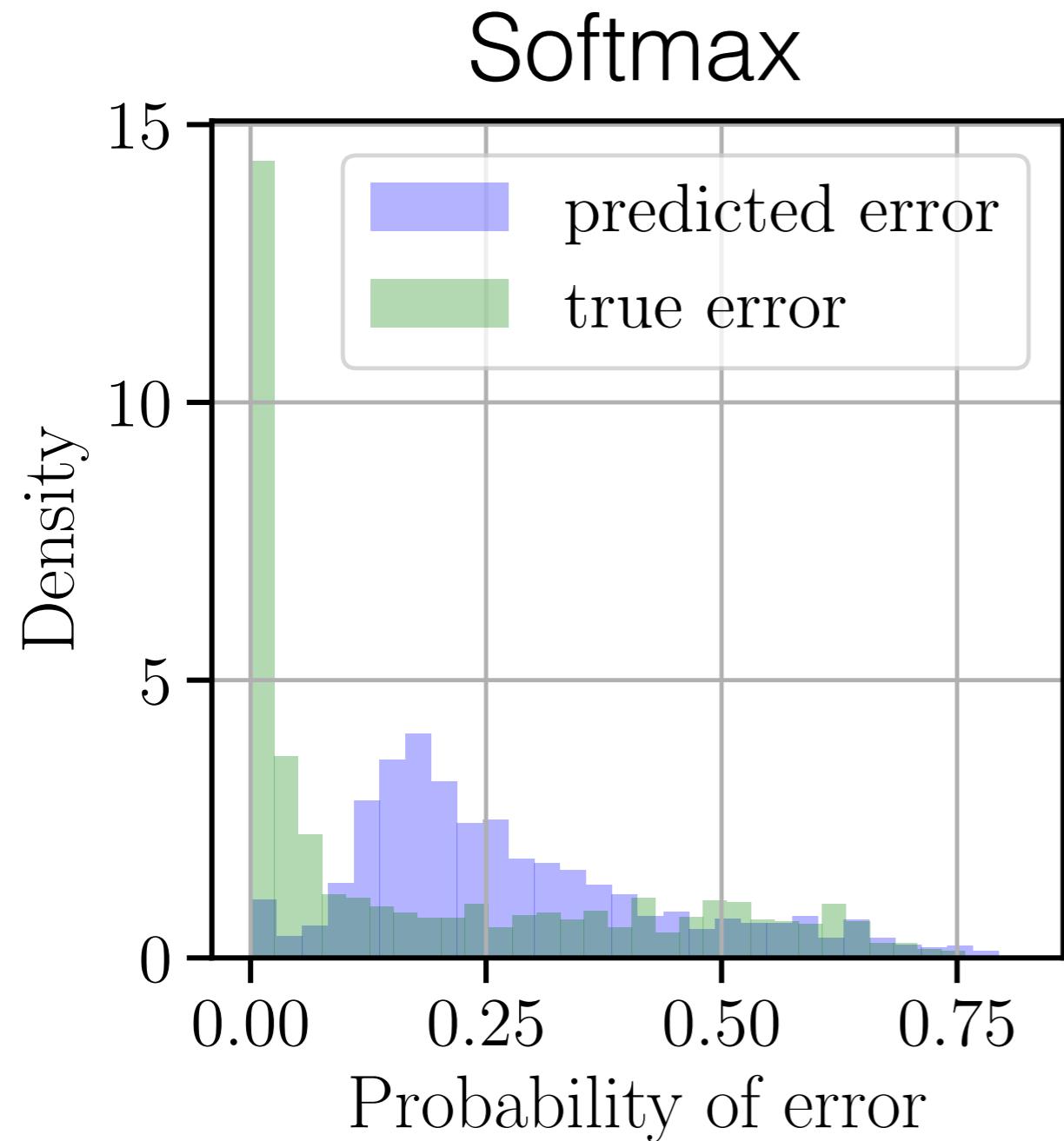
$$1 - P(m = y | x)$$

HAM10000: Skin Lesion Classification

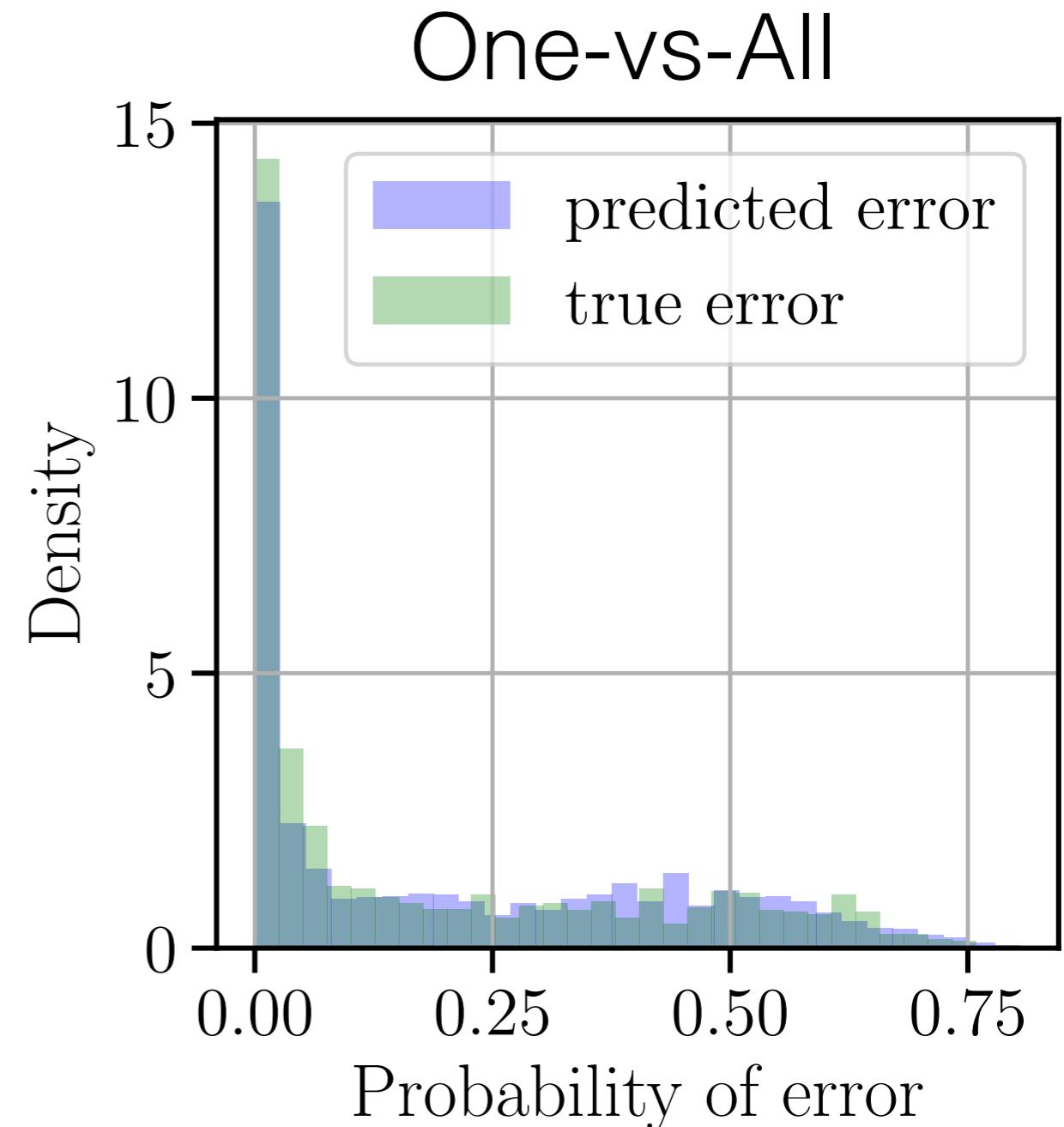


$$1 - P(m = y | x)$$

HAM10000: Skin Lesion Classification



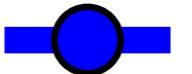
$$1 - P(m = y | x)$$



$$1 - P(m = y | x)$$

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

 one-vs-all (ours)

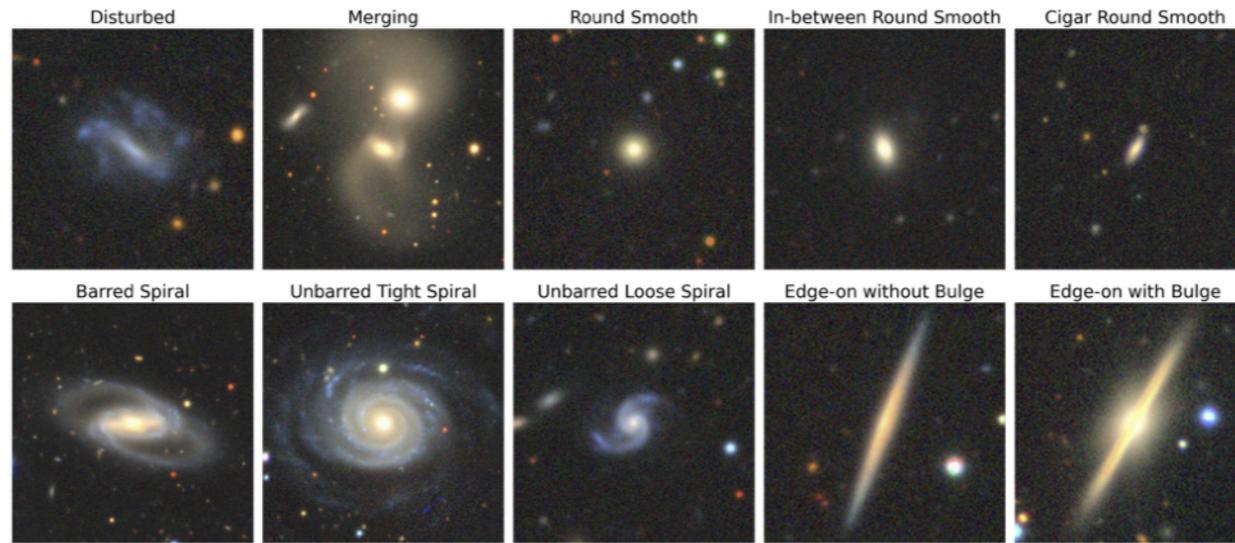
 softmax

 confidence

 score

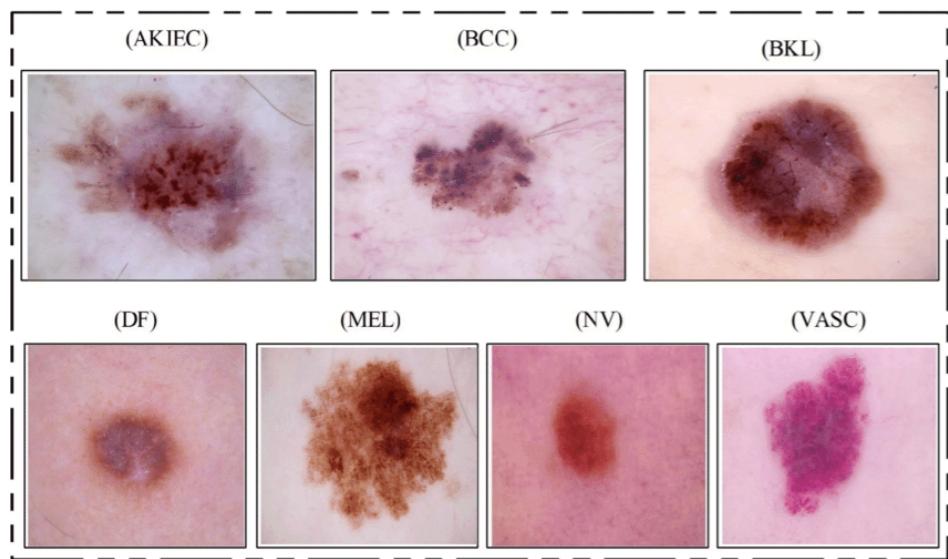
 differentiable triage

Galaxy Zoo



Galaxy10 DECs: Henry Leung/Jo Bovy 2021, Data: DECs/Galaxy Zoo

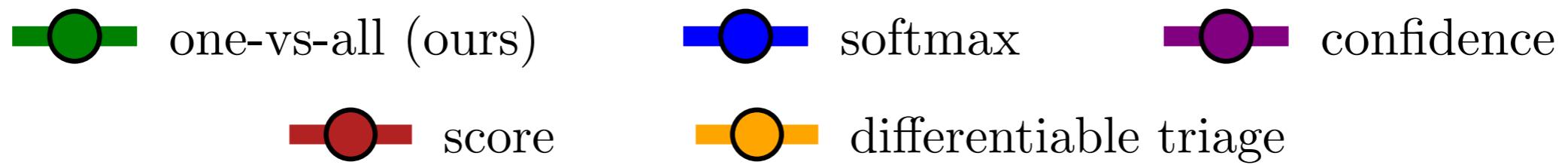
HAM10000



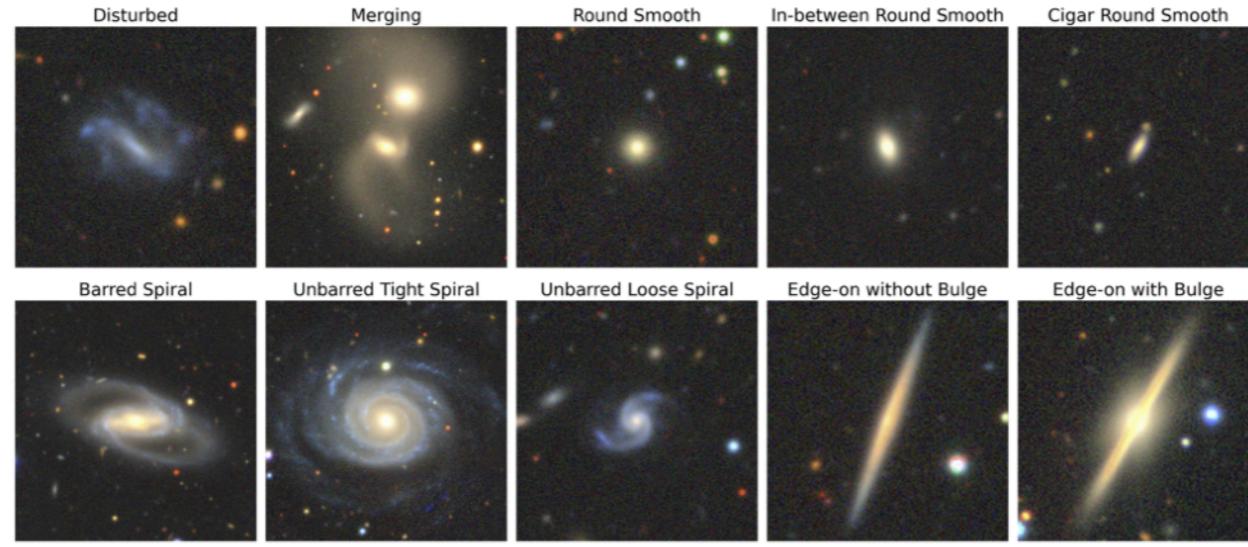
Hate Speech



[Davidson et al., ICWSM 2017]

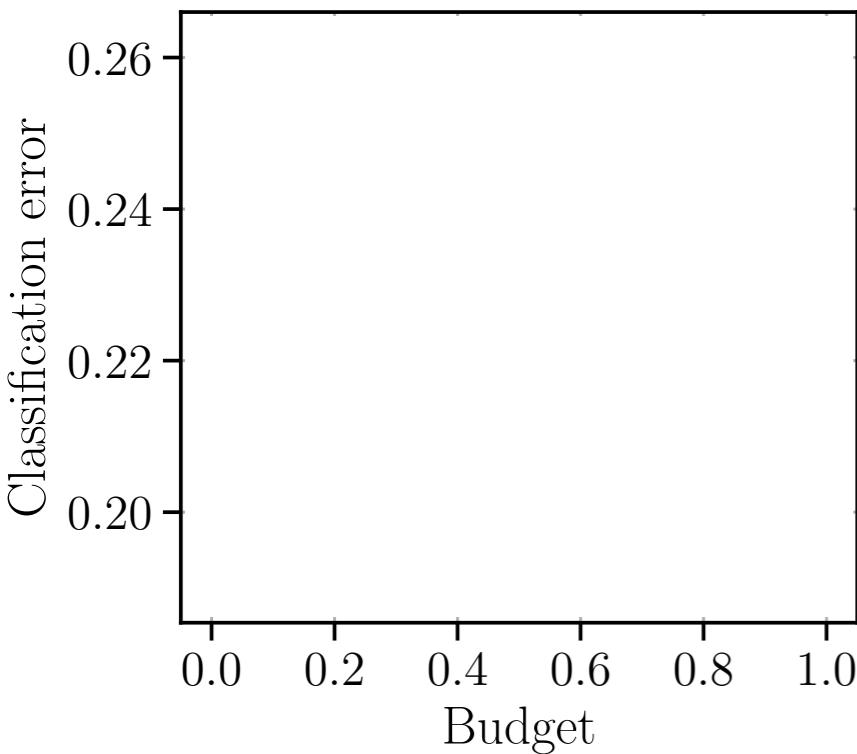


Galaxy Zoo



Galaxy10 DECalss: Henry Leung/Jo Bovy 2021, Data: DECalss/Galaxy Zoo

HAM10000



Hate Speech



[Davidson et al., ICWSM 2017]

 one-vs-all (ours)

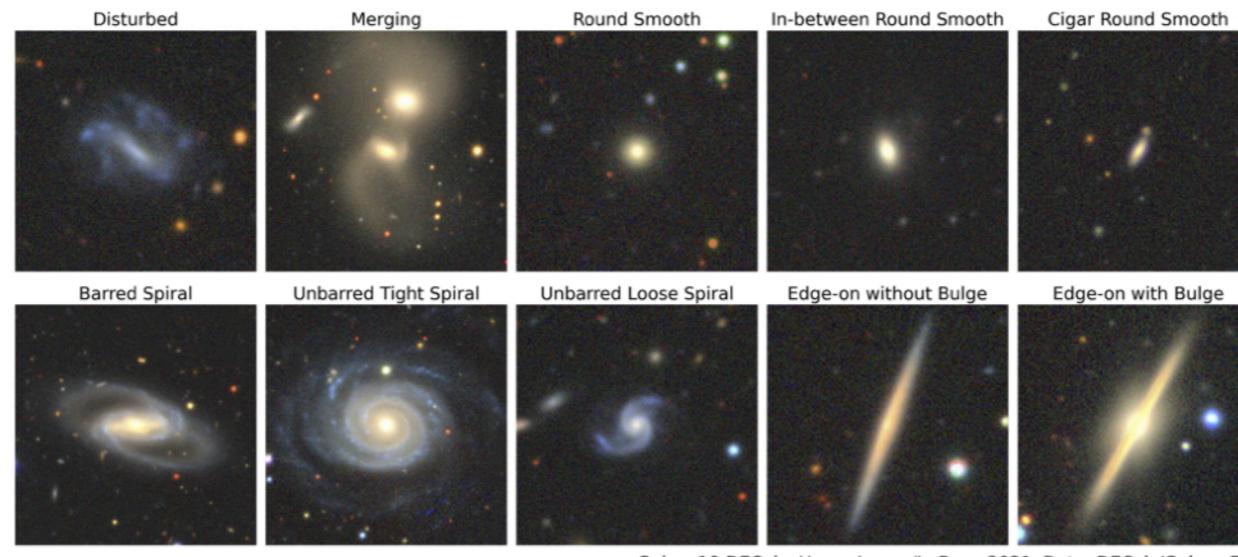
 softmax

 confidence

 score

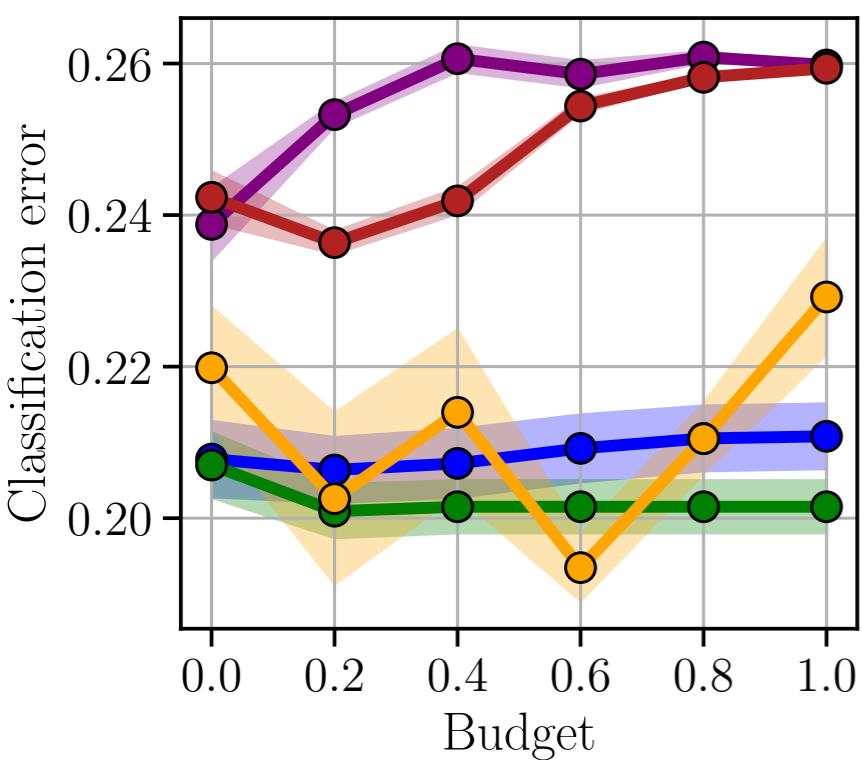
 differentiable triage

Galaxy Zoo



Galaxy10 DECs: Henry Leung/Jo Bovy 2021, Data: DECs/Galaxy Zoo

HAM10000



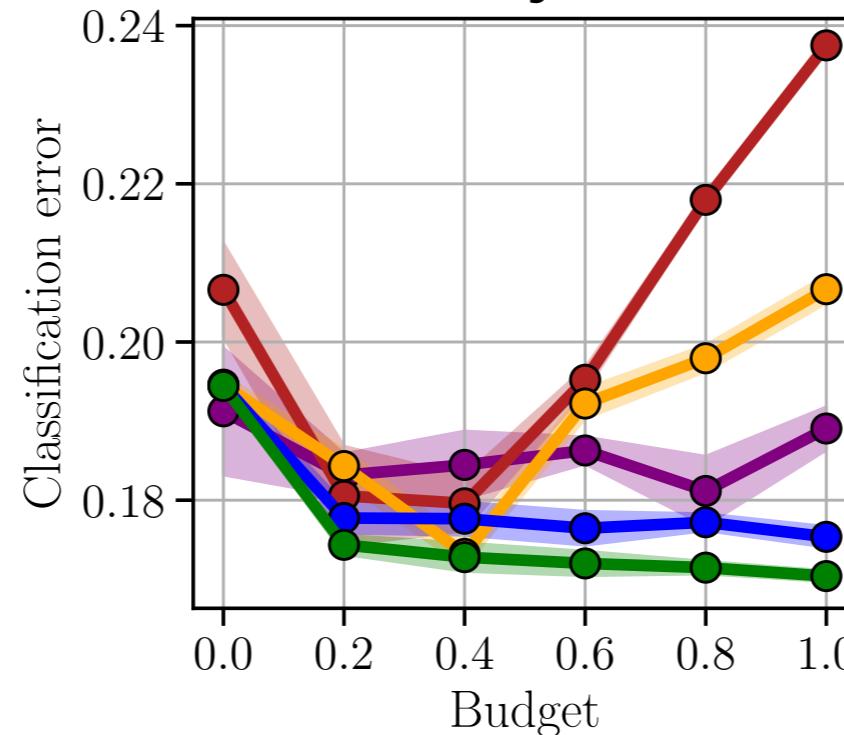
Hate Speech



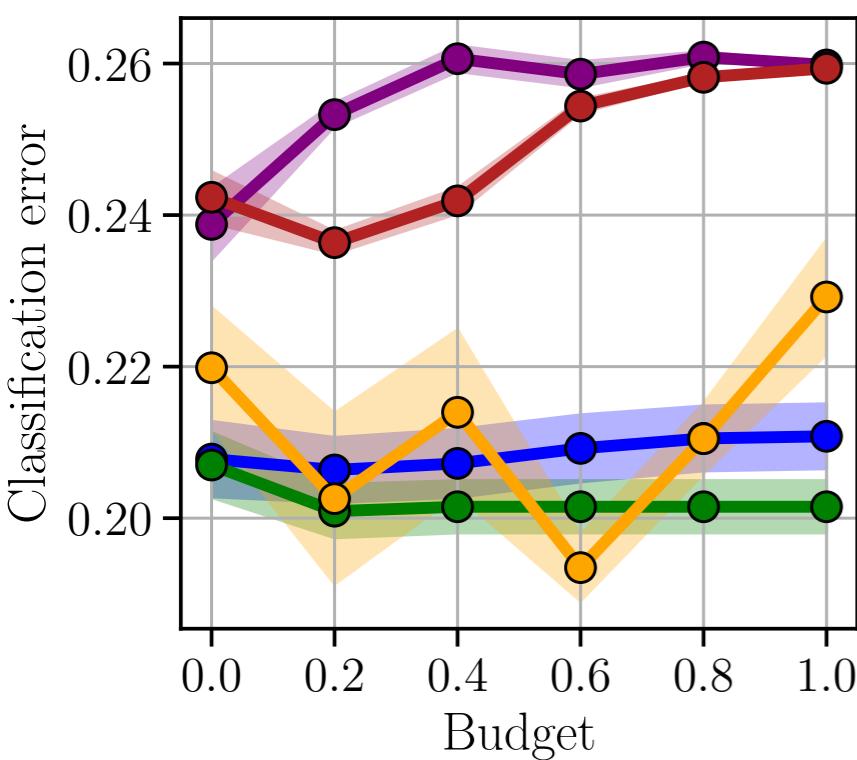
[Davidson et al., ICWSM 2017]

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

Galaxy Zoo



HAM10000



Hate Speech



[Davidson et al., ICWSM 2017]

 one-vs-all (ours)

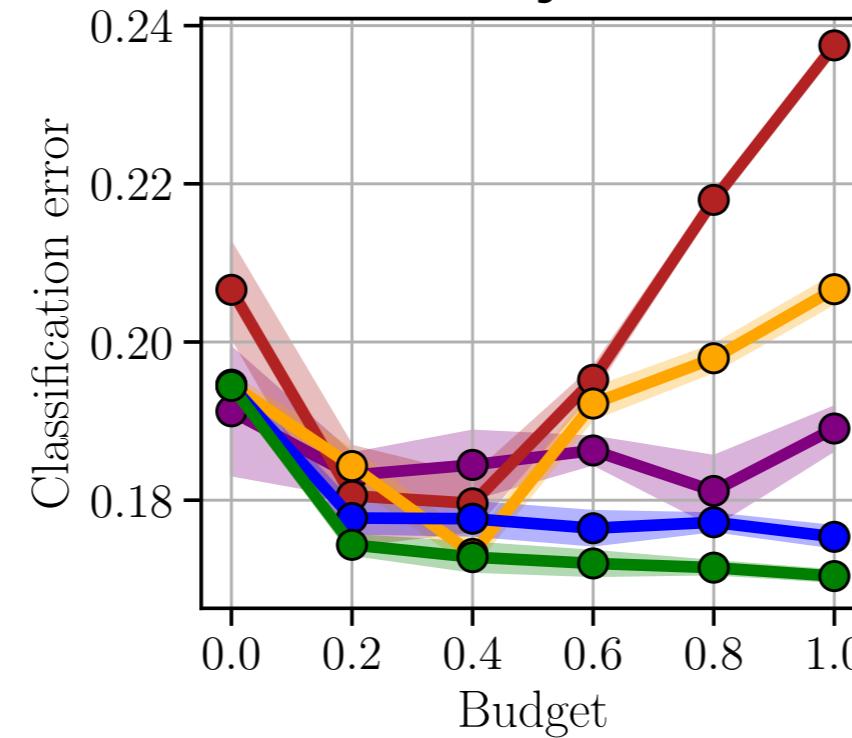
 softmax

 confidence

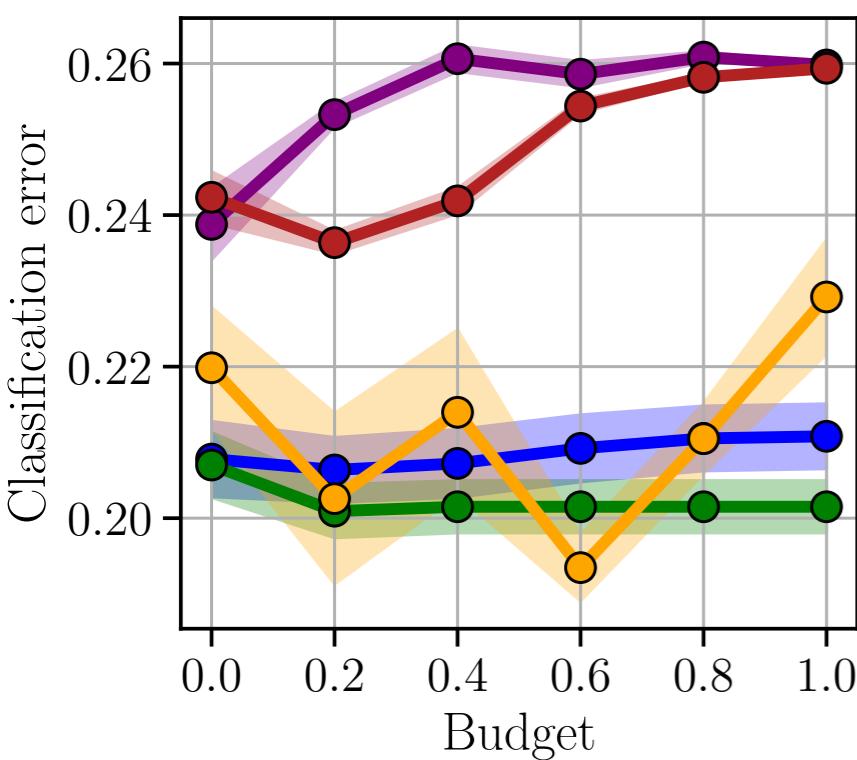
 score

 differentiable triage

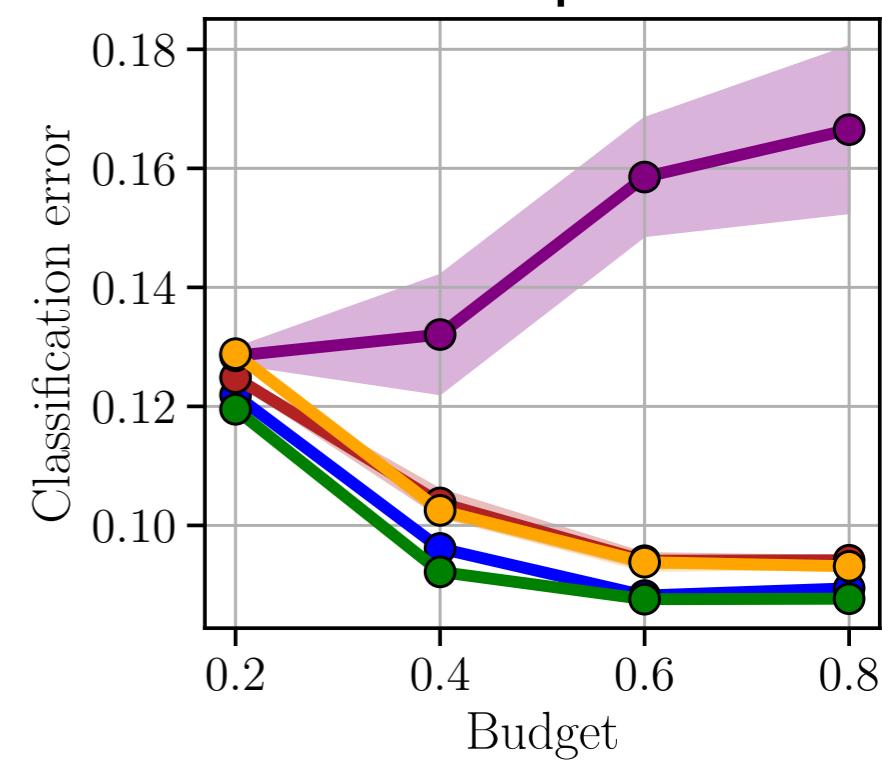
Galaxy Zoo

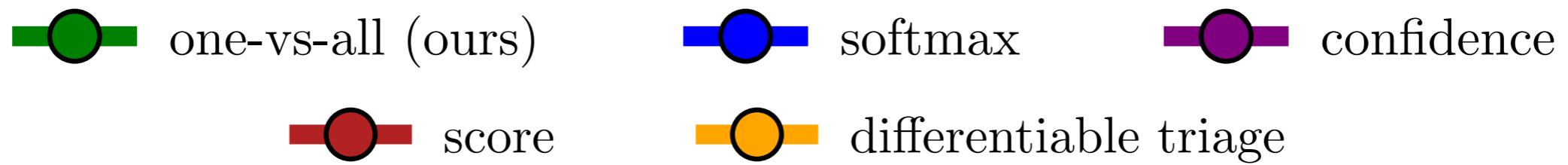


HAM10000

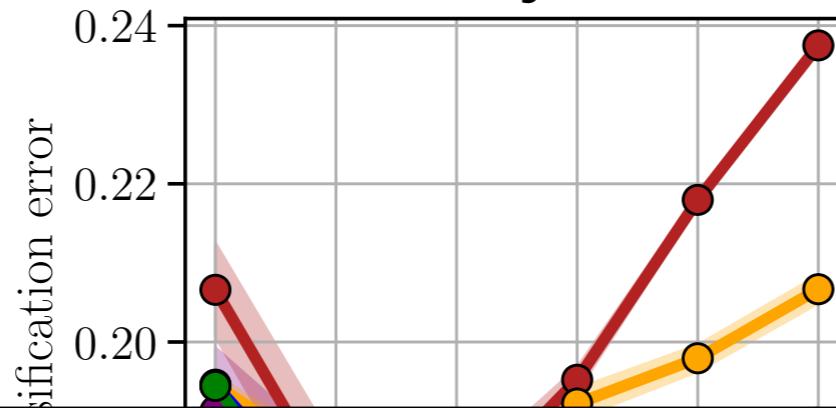


Hate Speech

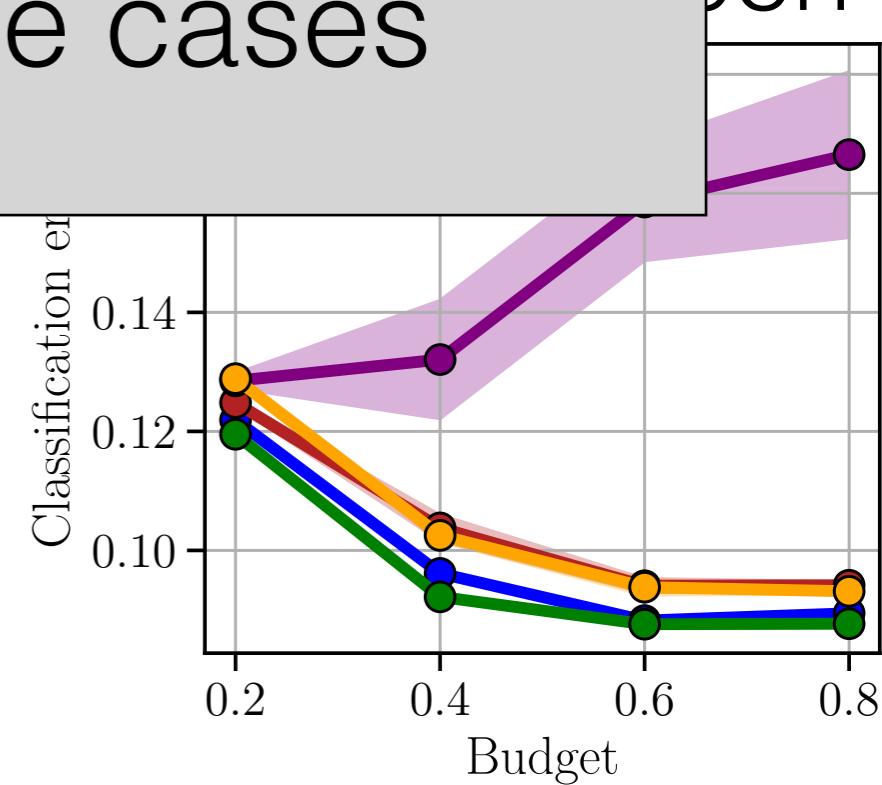
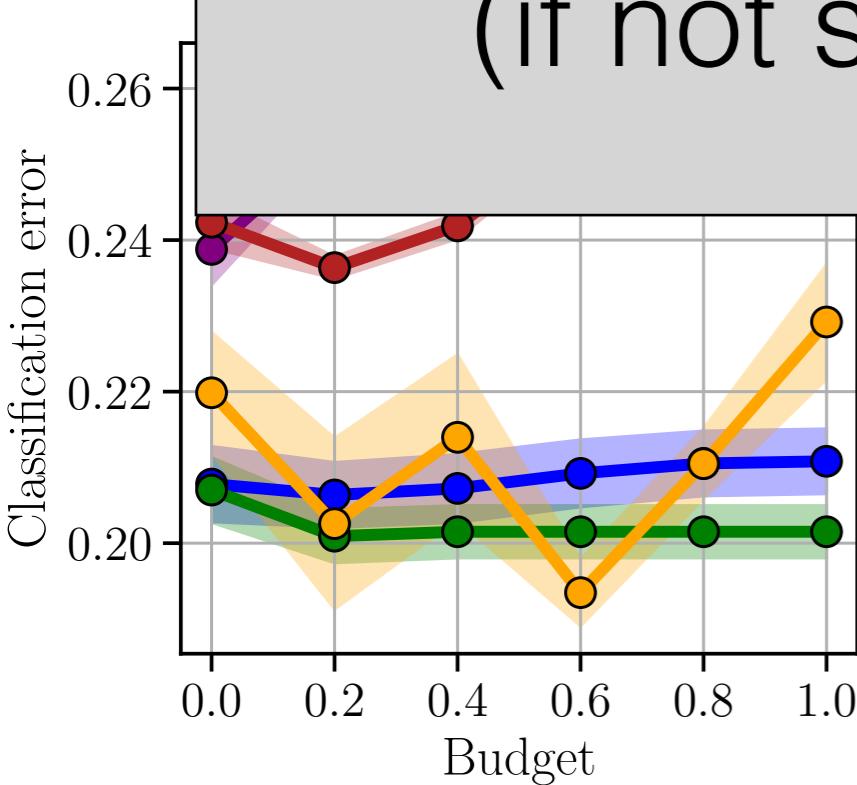




Galaxy Zoo

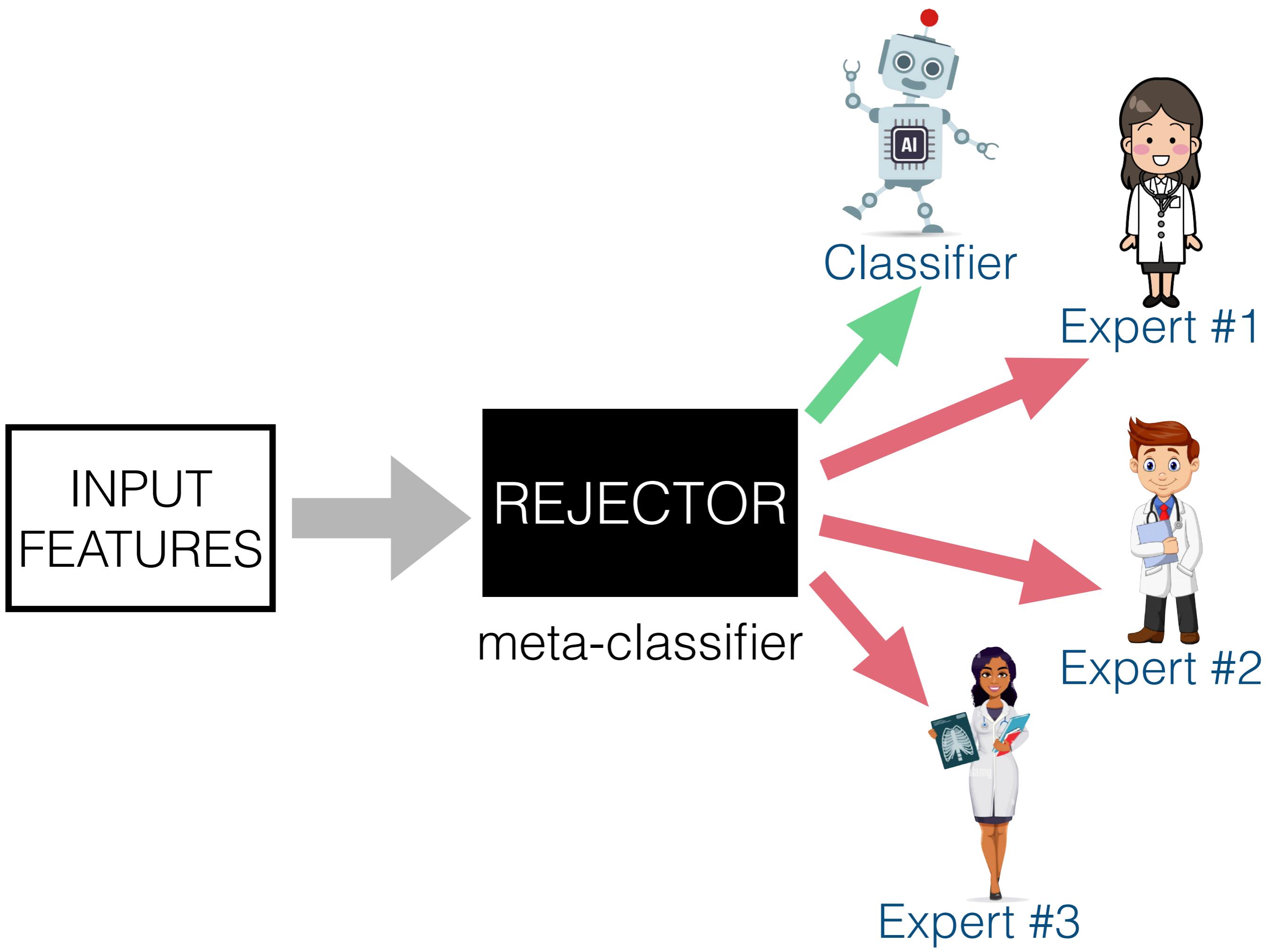


One-vs-all models are competitive
(if not superior) in all three cases

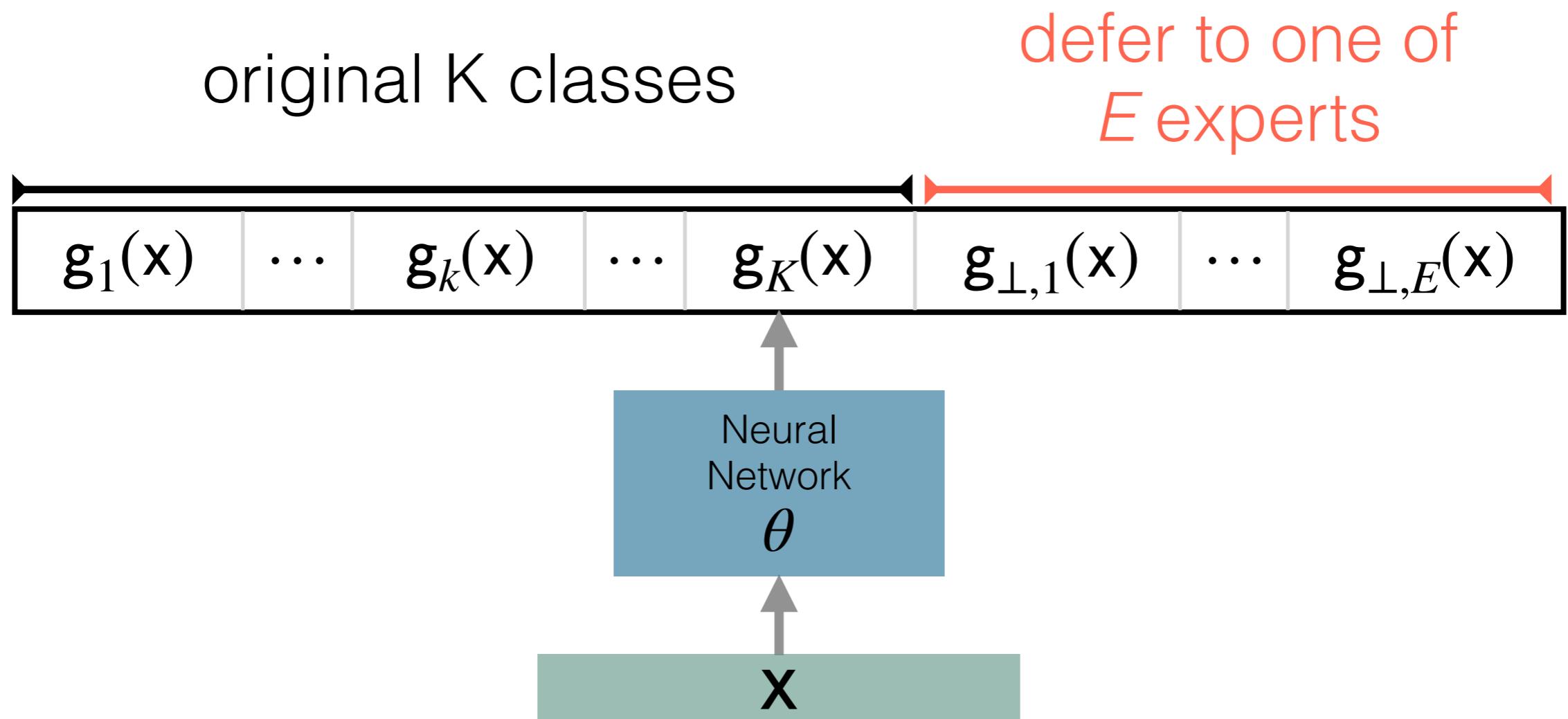


Extension

Learning to Defer
to Multiple Experts



Learning-to-Defer to Multiple Experts



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

$$-\sum_n \left(\log p_{y_n}(x_n) + \sum_e \mathbb{I}[y_n = m_{n,e}] \log p_{\perp,e}(x_n) \right)$$

classifier loss



rejector loss

only if e -th expert is correct

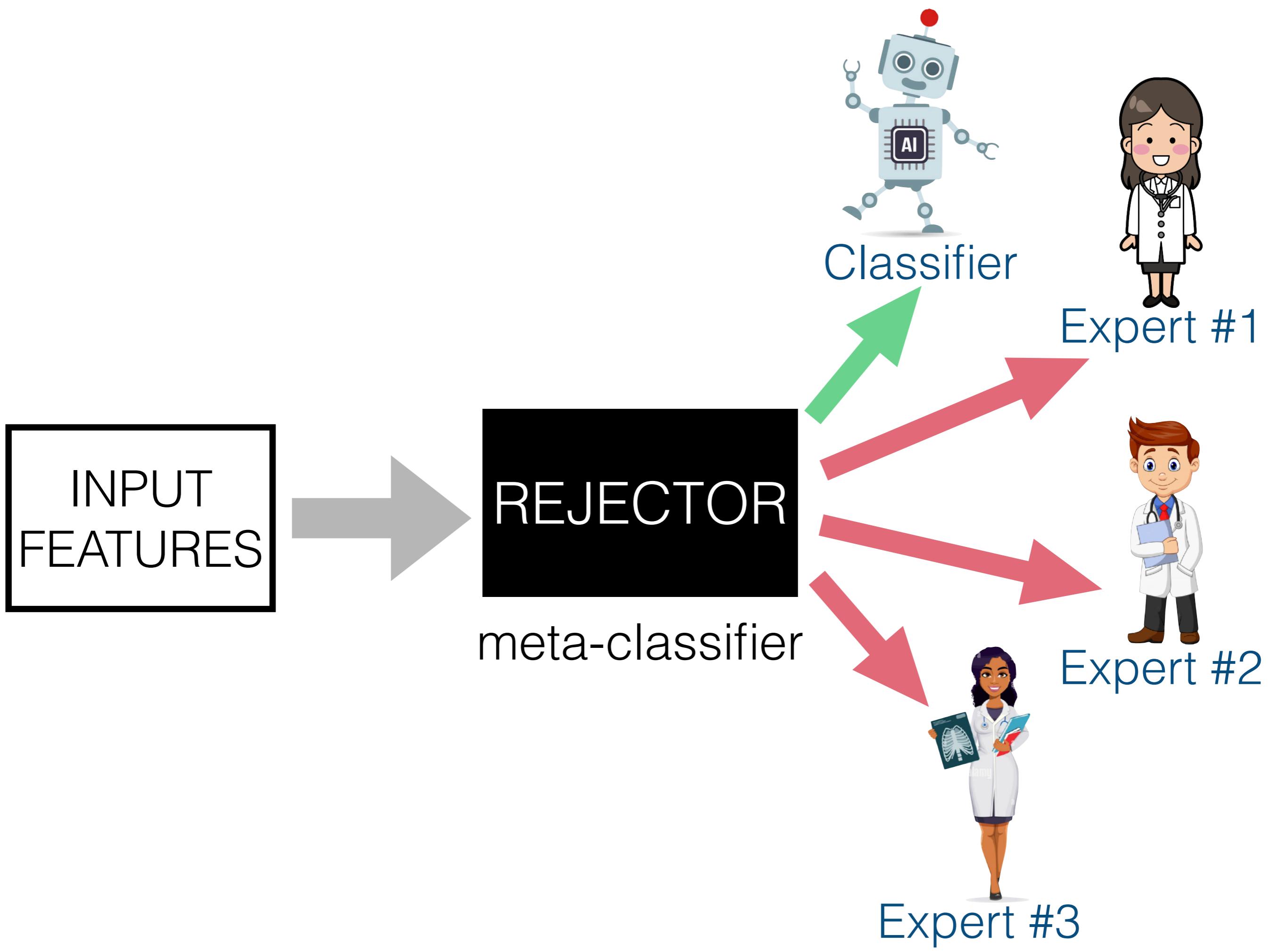
Calibration: Is the system a good forecaster?

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

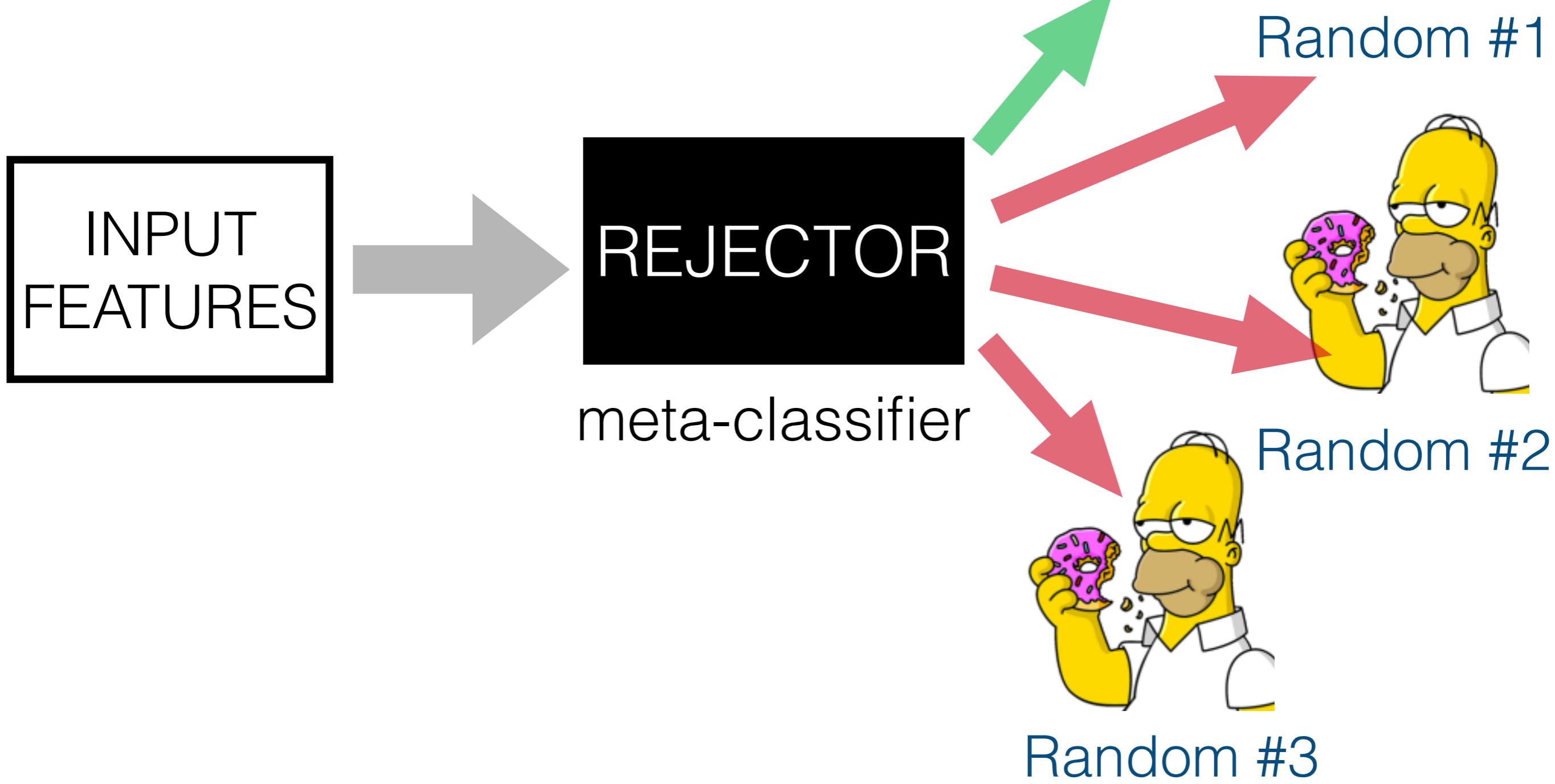
Does the rejector correctly estimate the e -th expert's chance of being correct?

Calibration: Is the system a good forecaster?

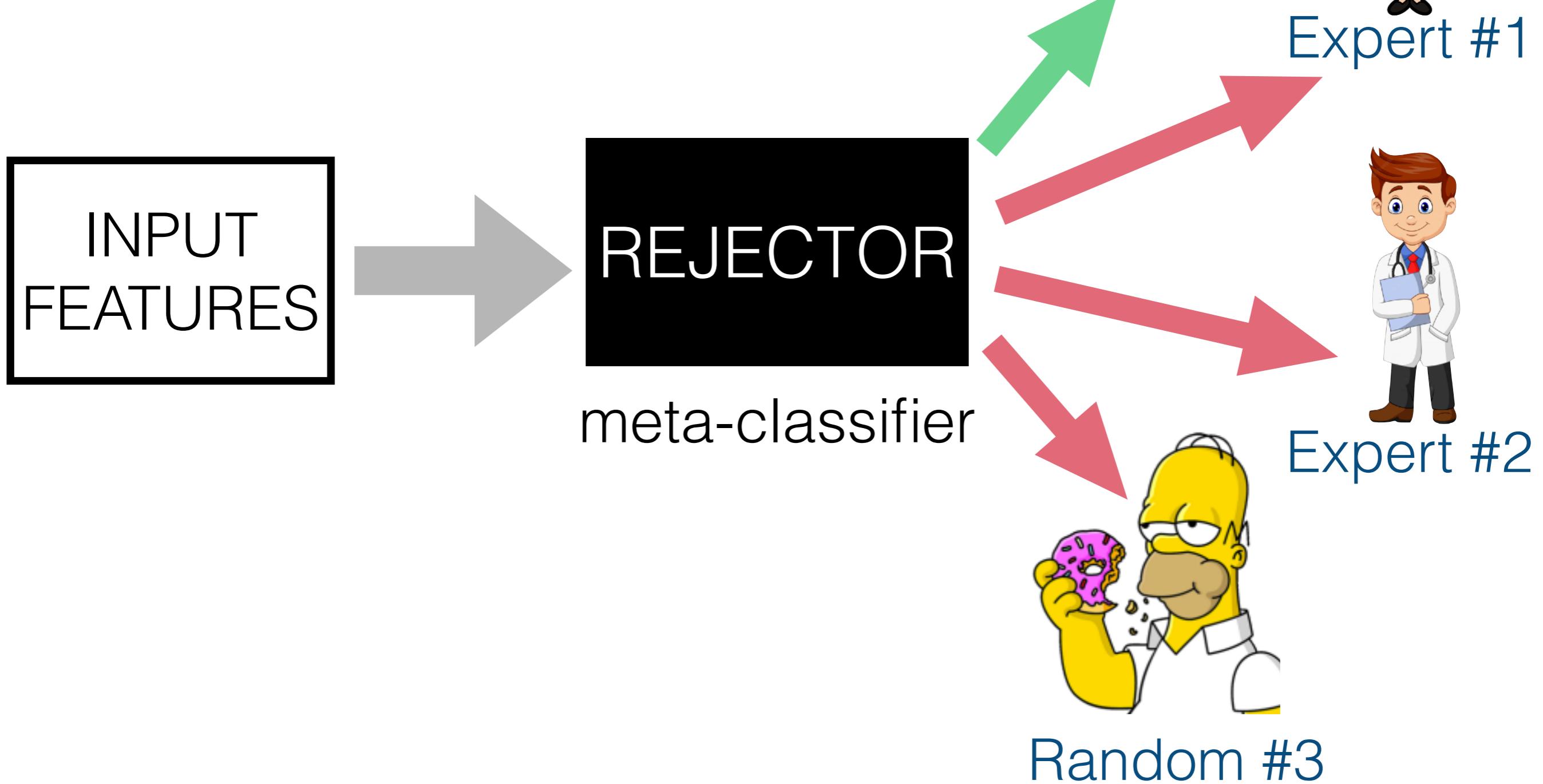
$$\mathbb{P}(m_e = y | x) = \frac{p_{\perp,e}^*(x)}{1 - \sum_{j=1}^E p_{\perp,j}^*(x)}$$



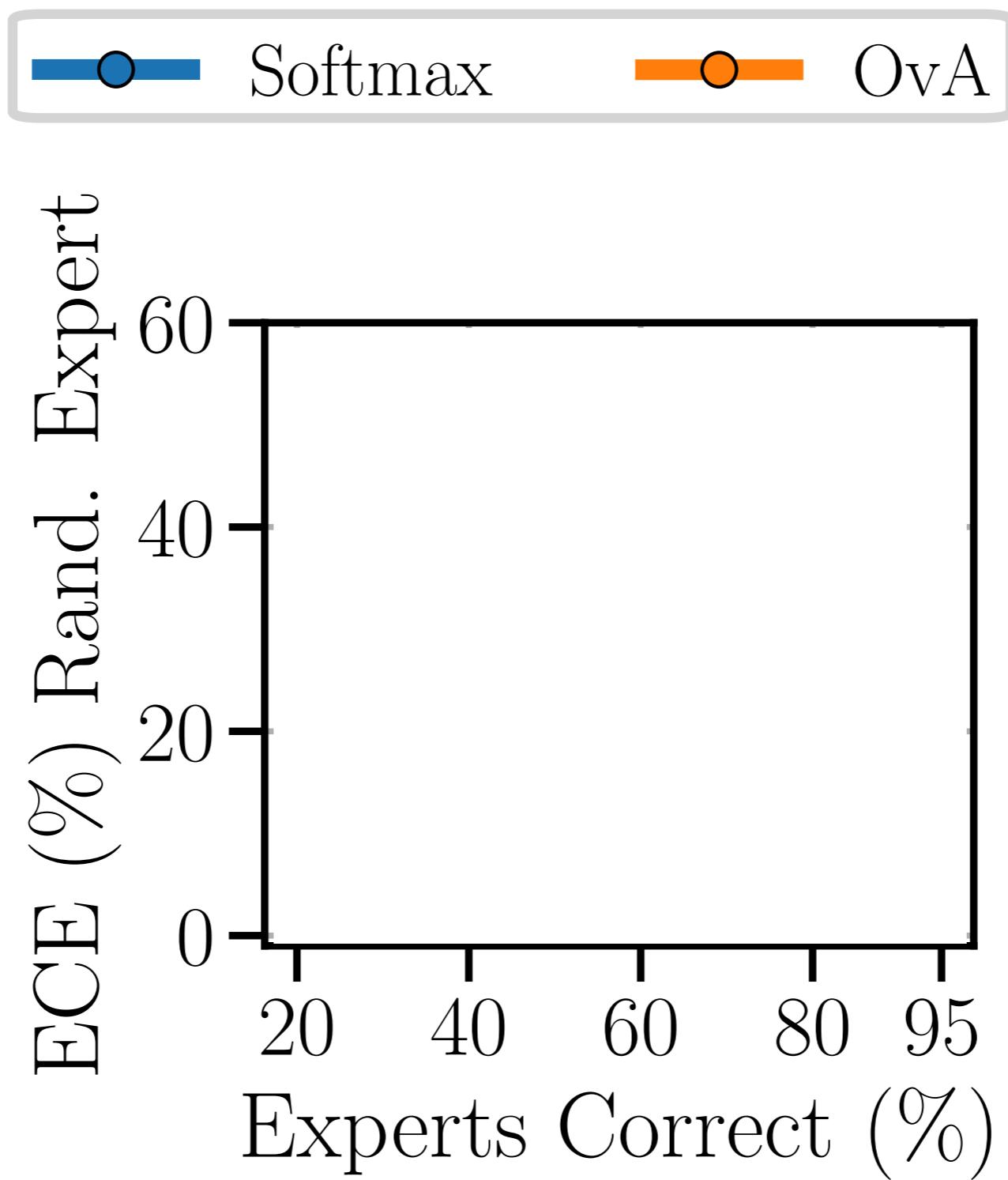
Training: random experts



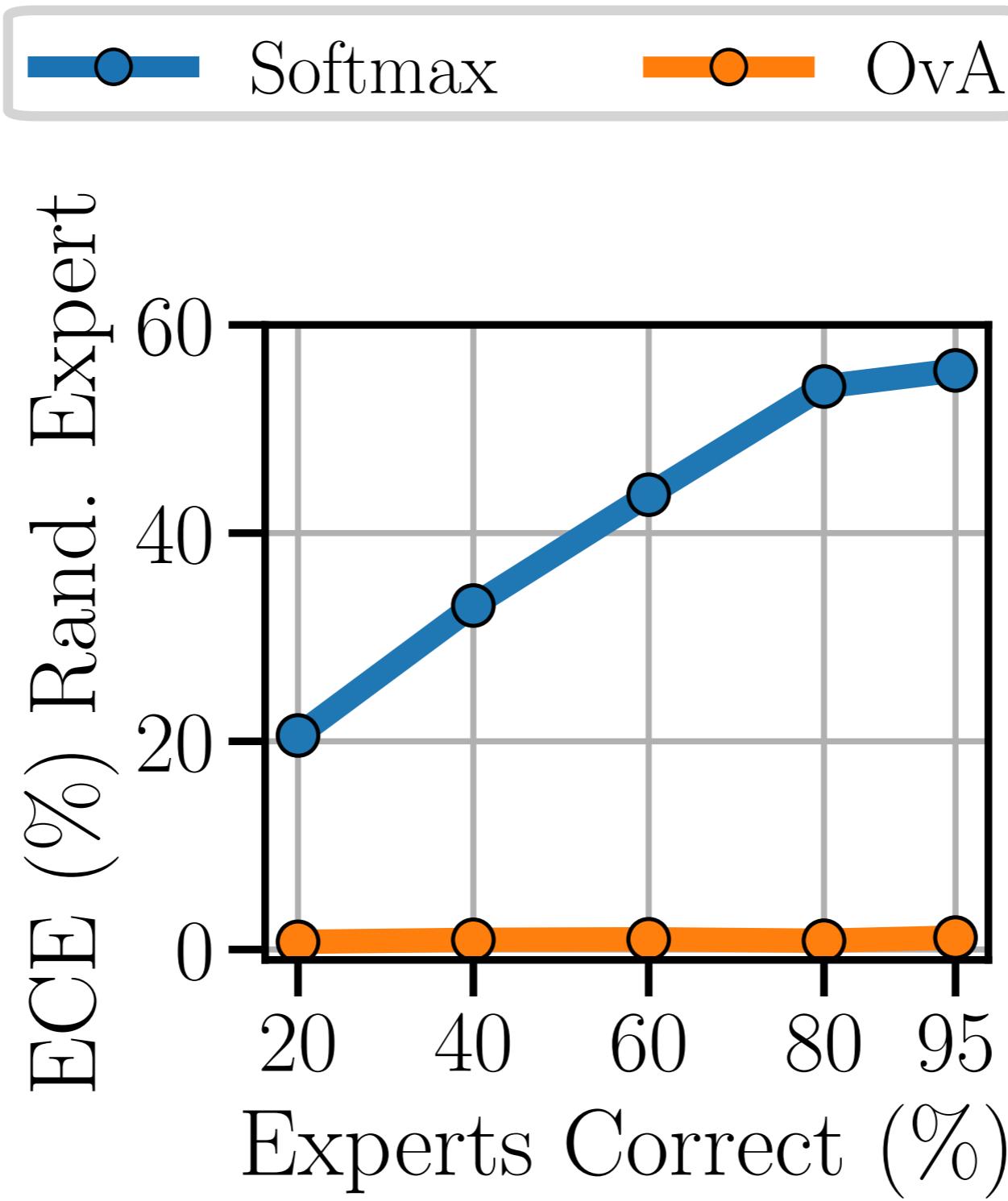
Testing: one random



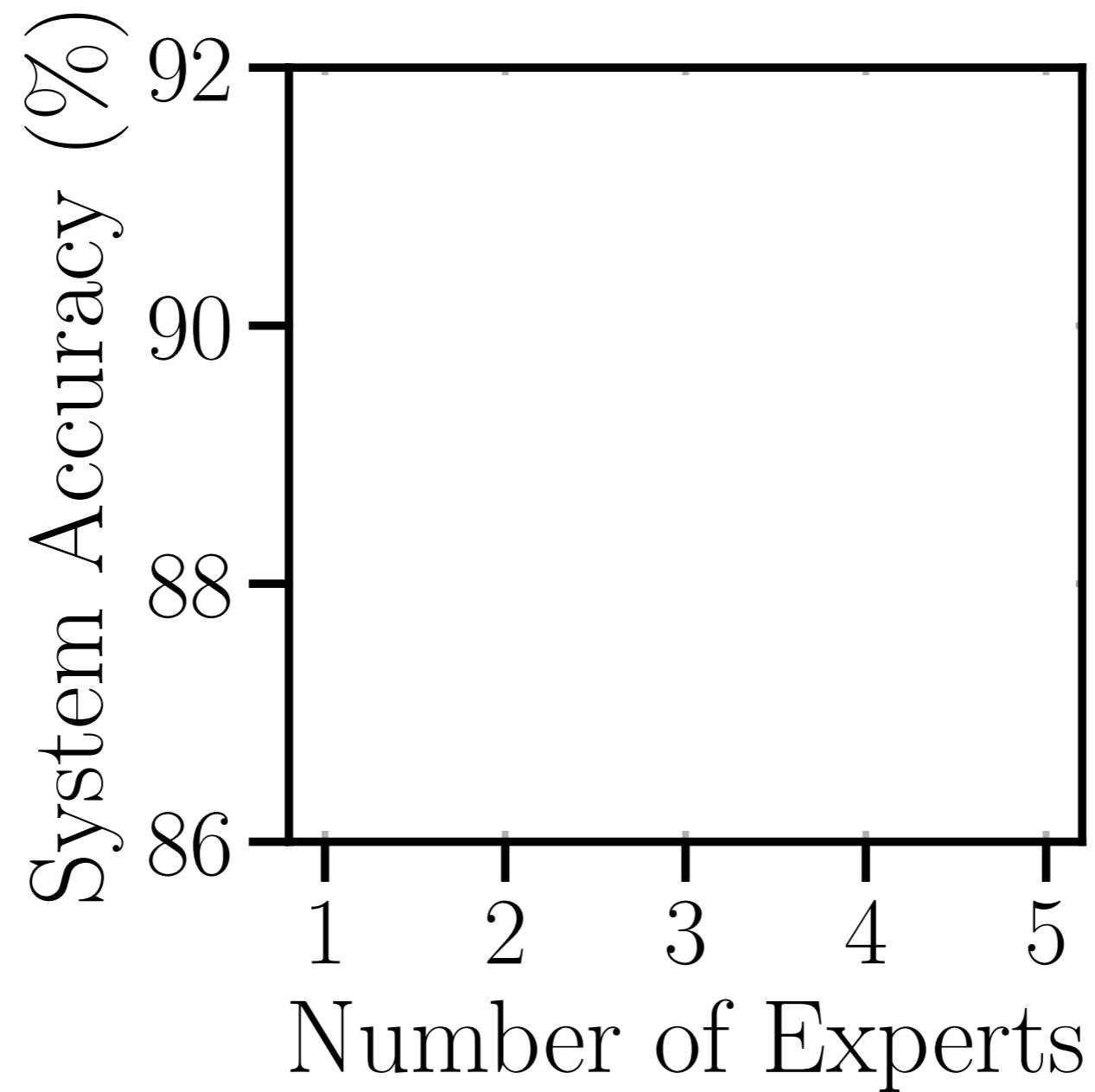
CIFAR-10: Expert Dependence



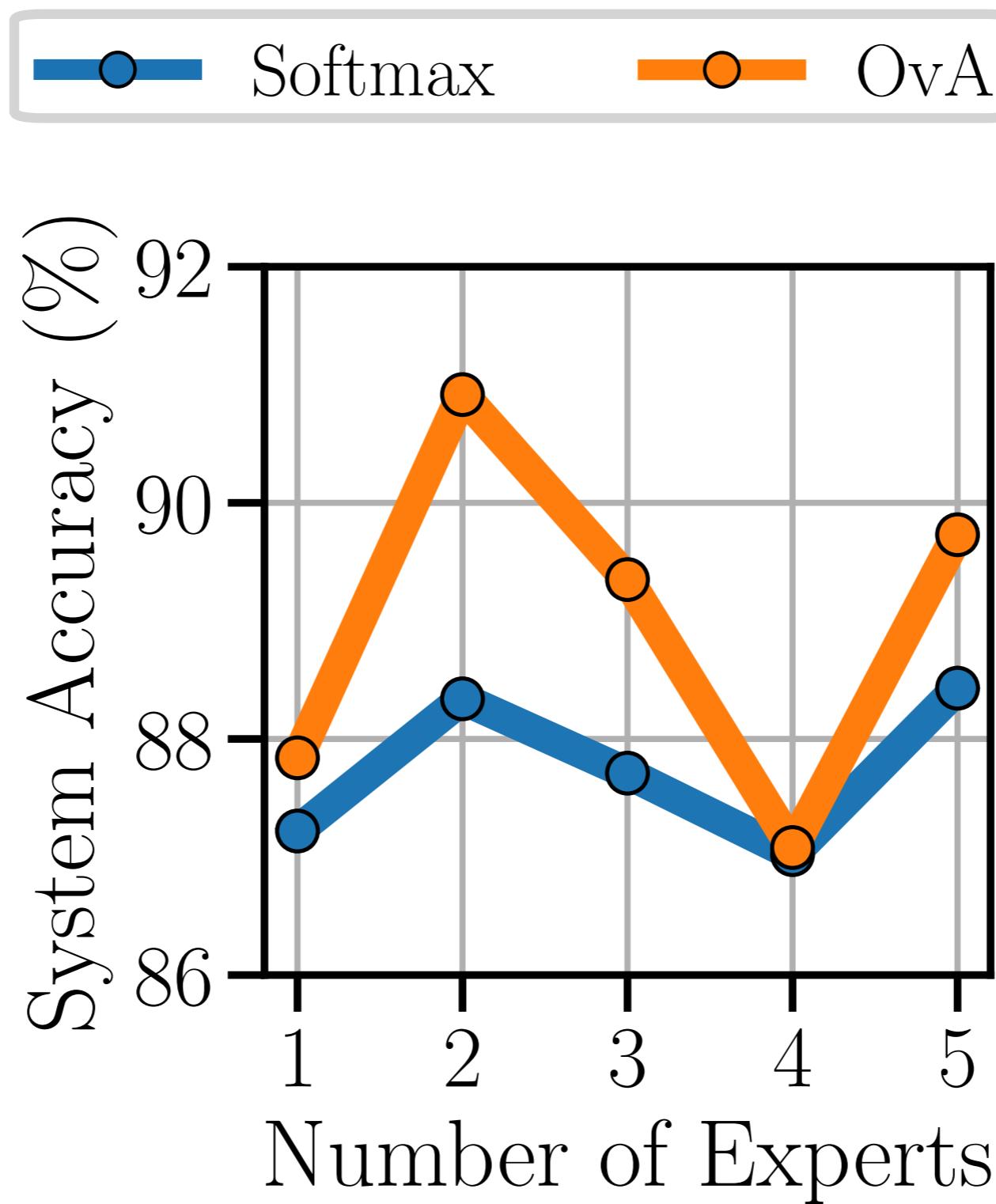
CIFAR-10: Expert Dependence



CIFAR-10: Specialized Experts



CIFAR-10: Specialized Experts



Summary

- ⊗ Softmax-based losses can produce invalid estimators for expert correctness.
- ⊗ One-vs-all formulation retains theoretical and practical benefits while having calibrated estimates of expert correctness.
- ⊗ Problem for softmax gets worse when there are multiple experts.

Future Work

- ⊗ Heavy dependence on supervised data.

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


- ⊗ Calibration will be essential for active learning, experimental design, missing data...

Paper



Code



Funding



Thank you. Questions?