

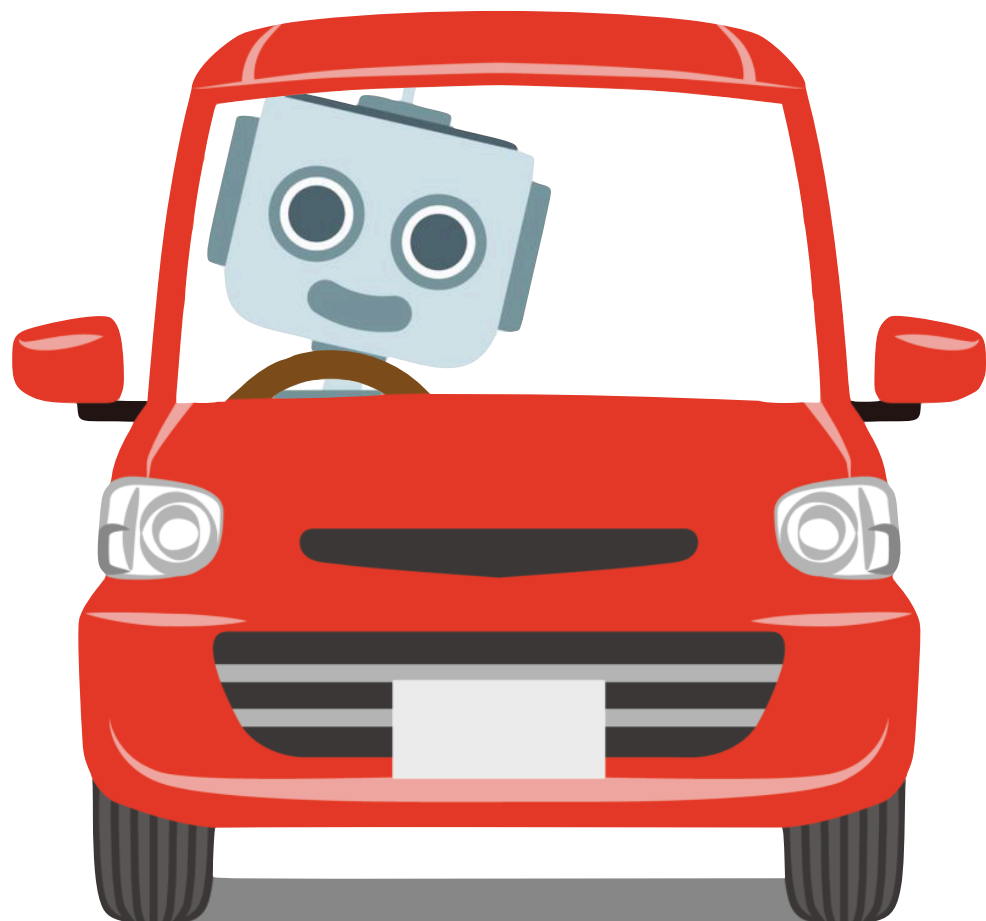
# Learning to Defer to One, Multiple, or a Population of Expert(s)

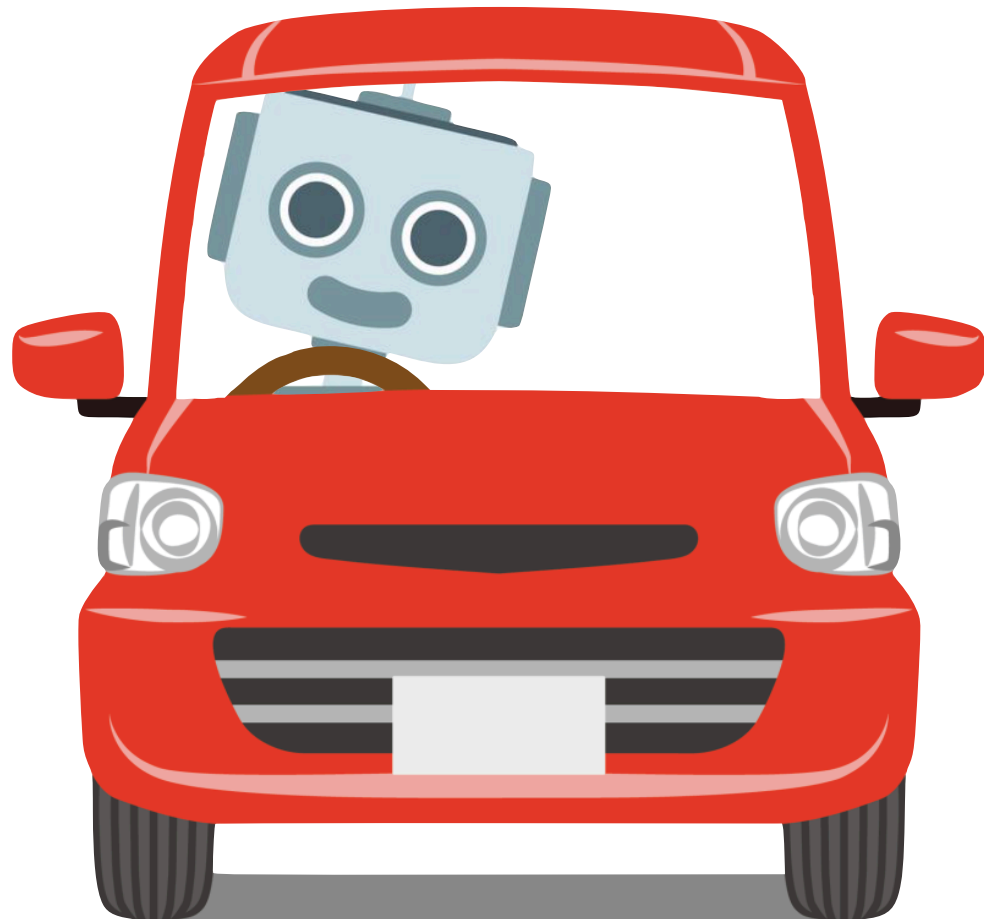
Eric Nalisnick

Johns Hopkins University

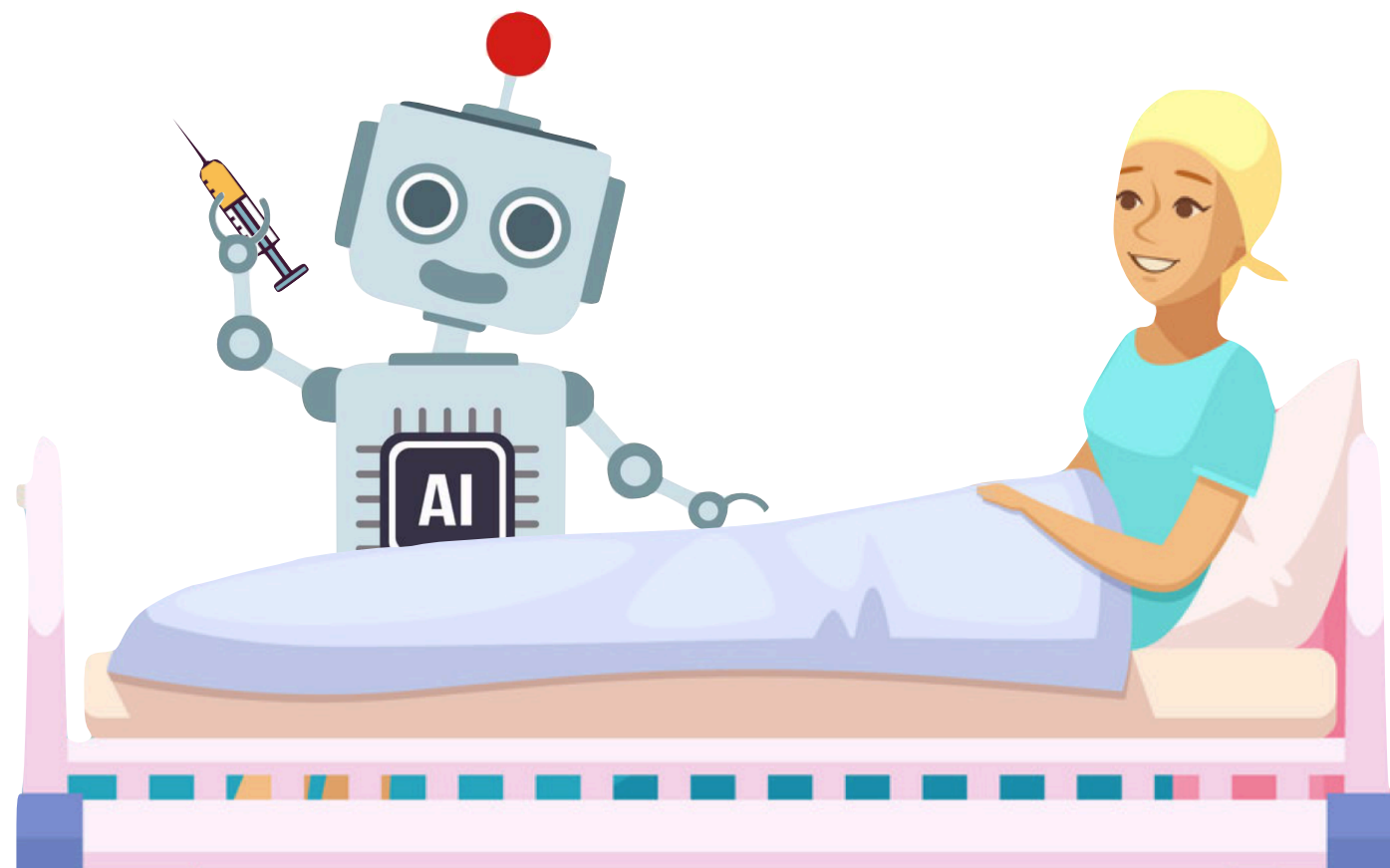
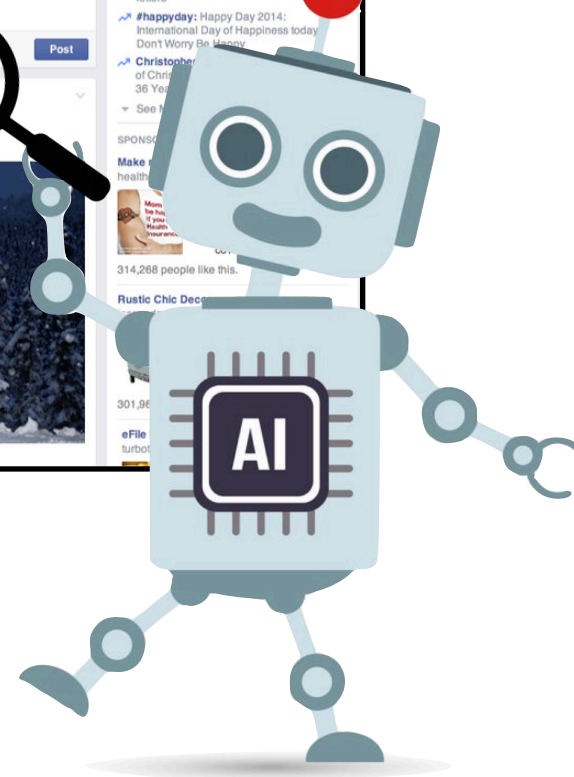
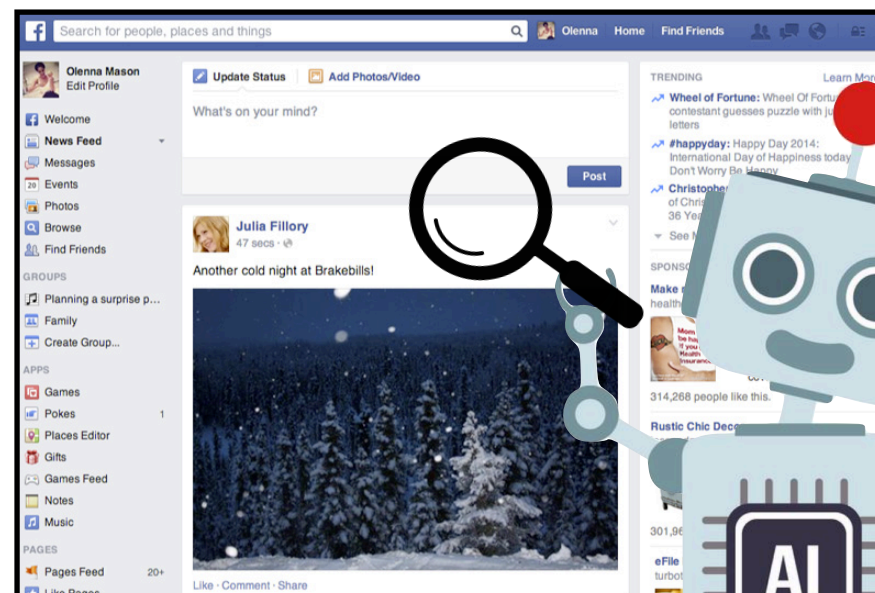
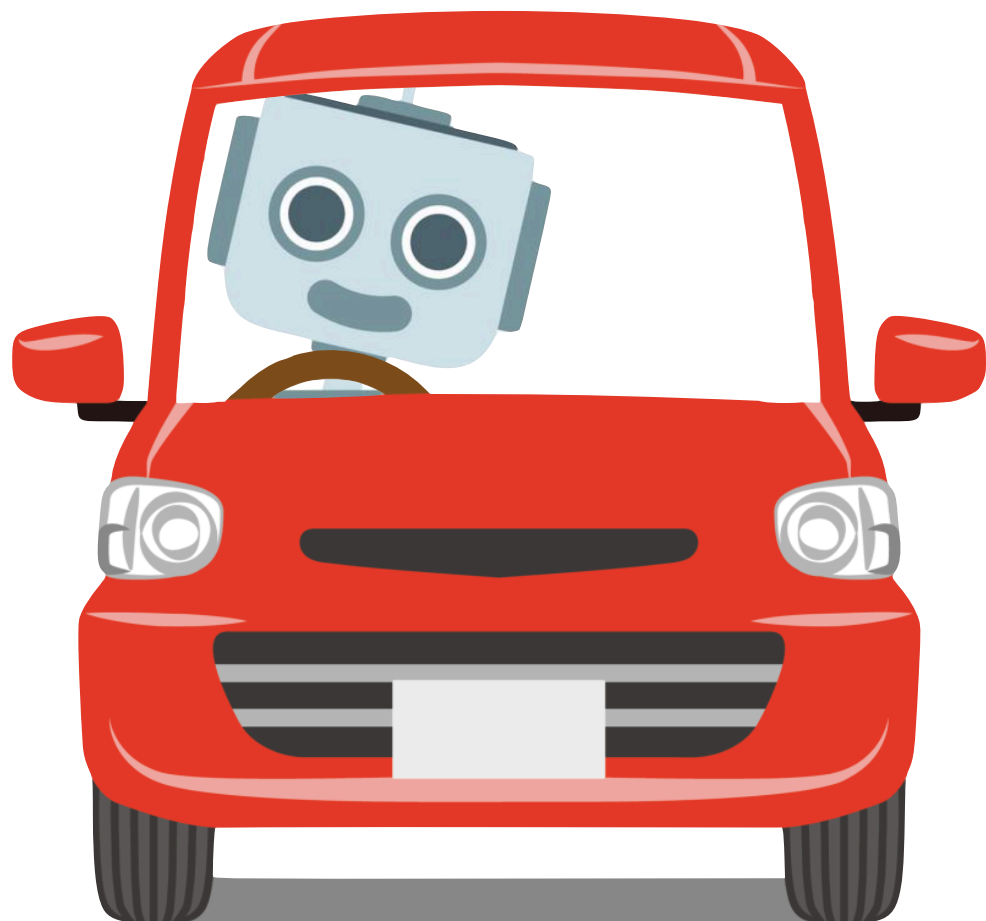


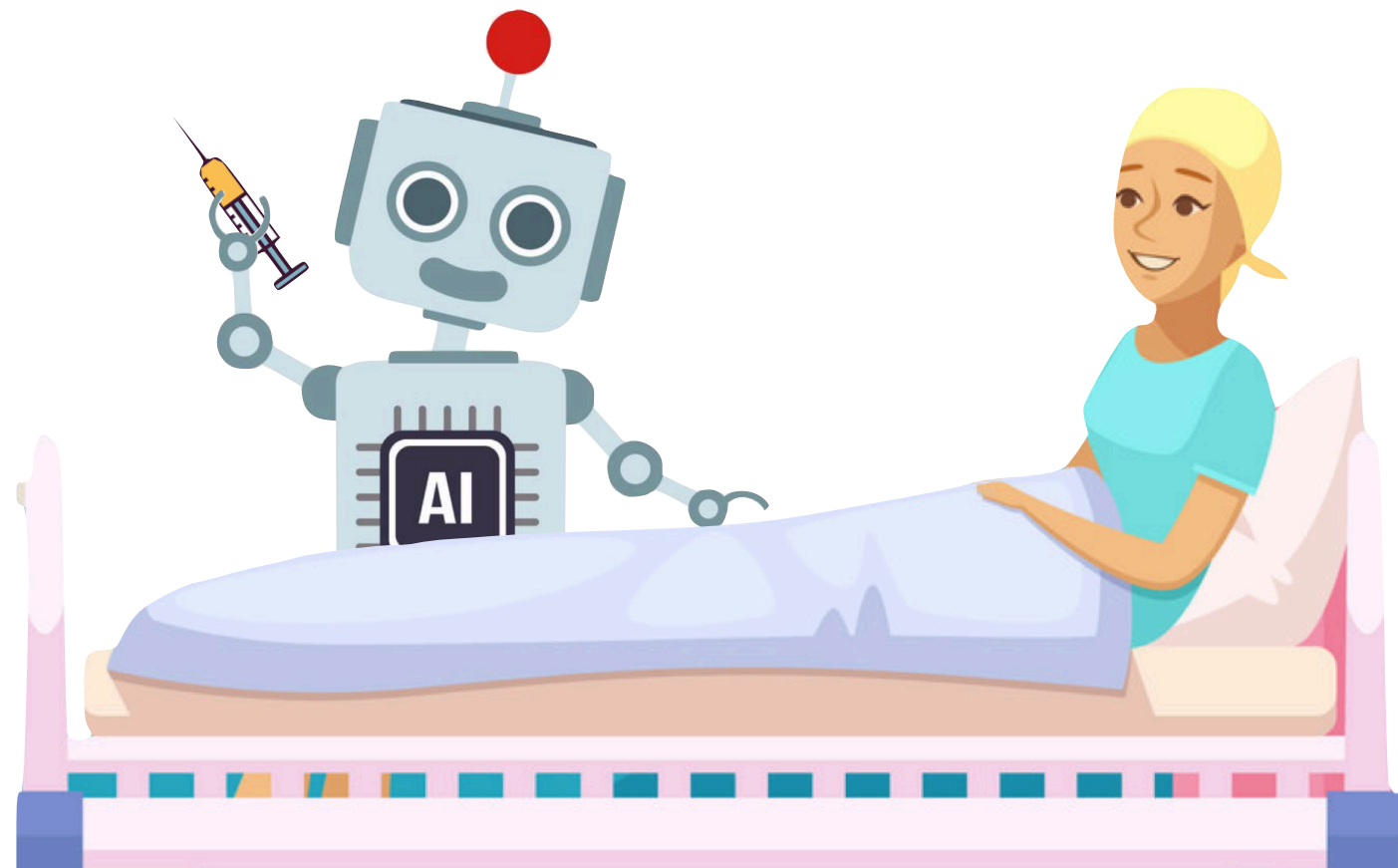
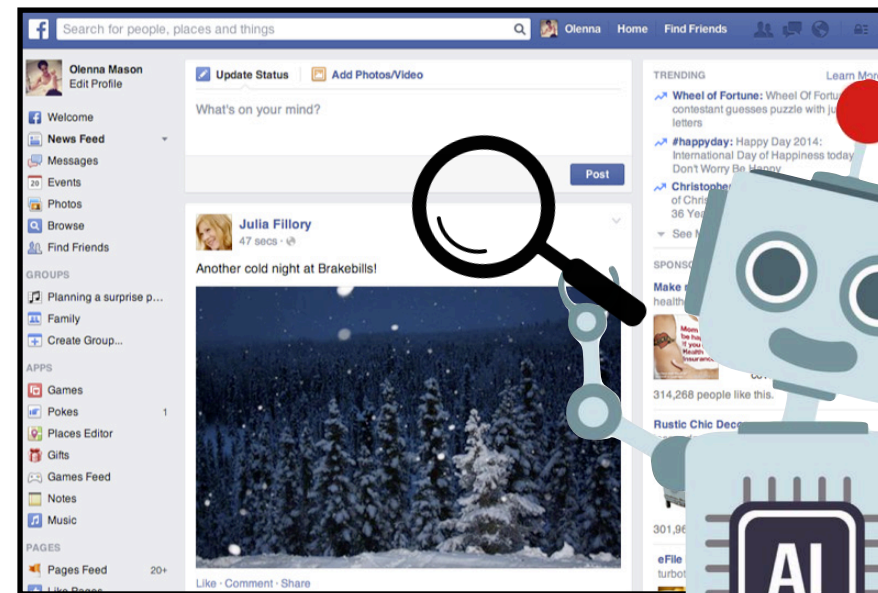
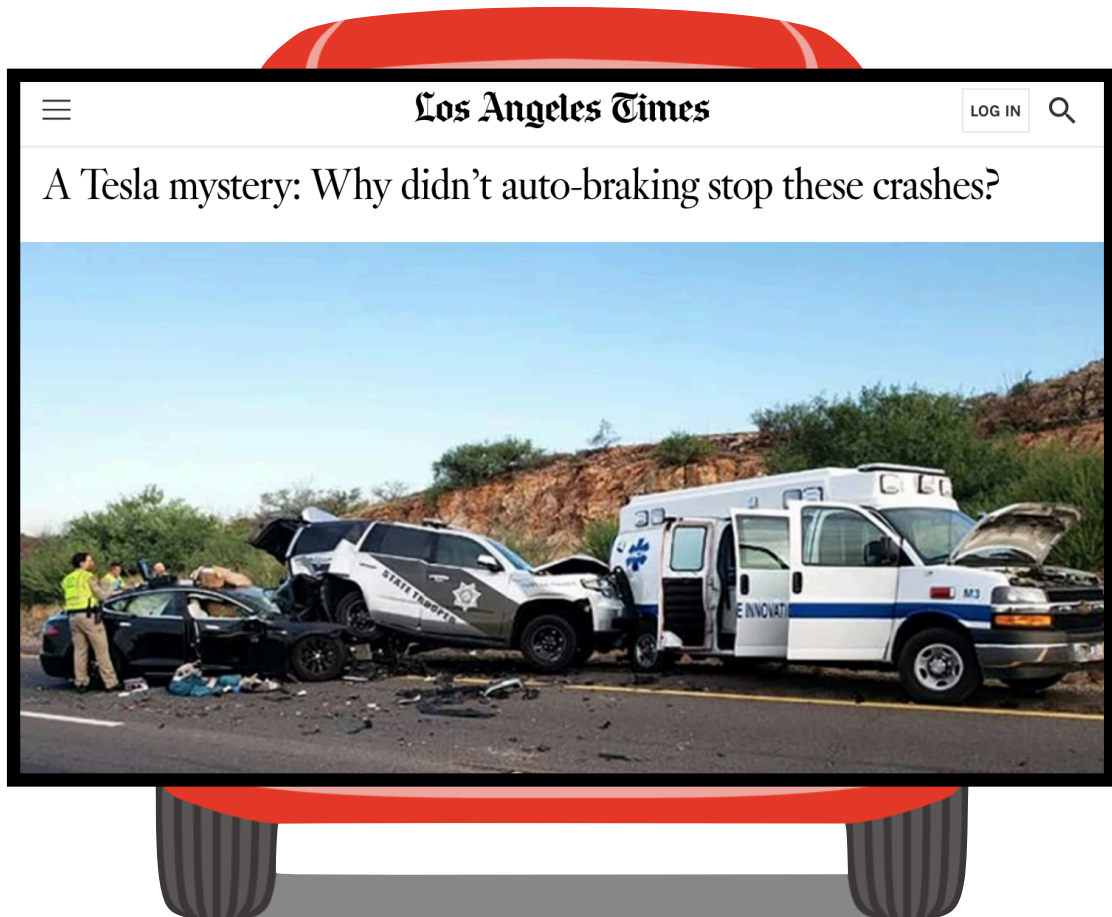




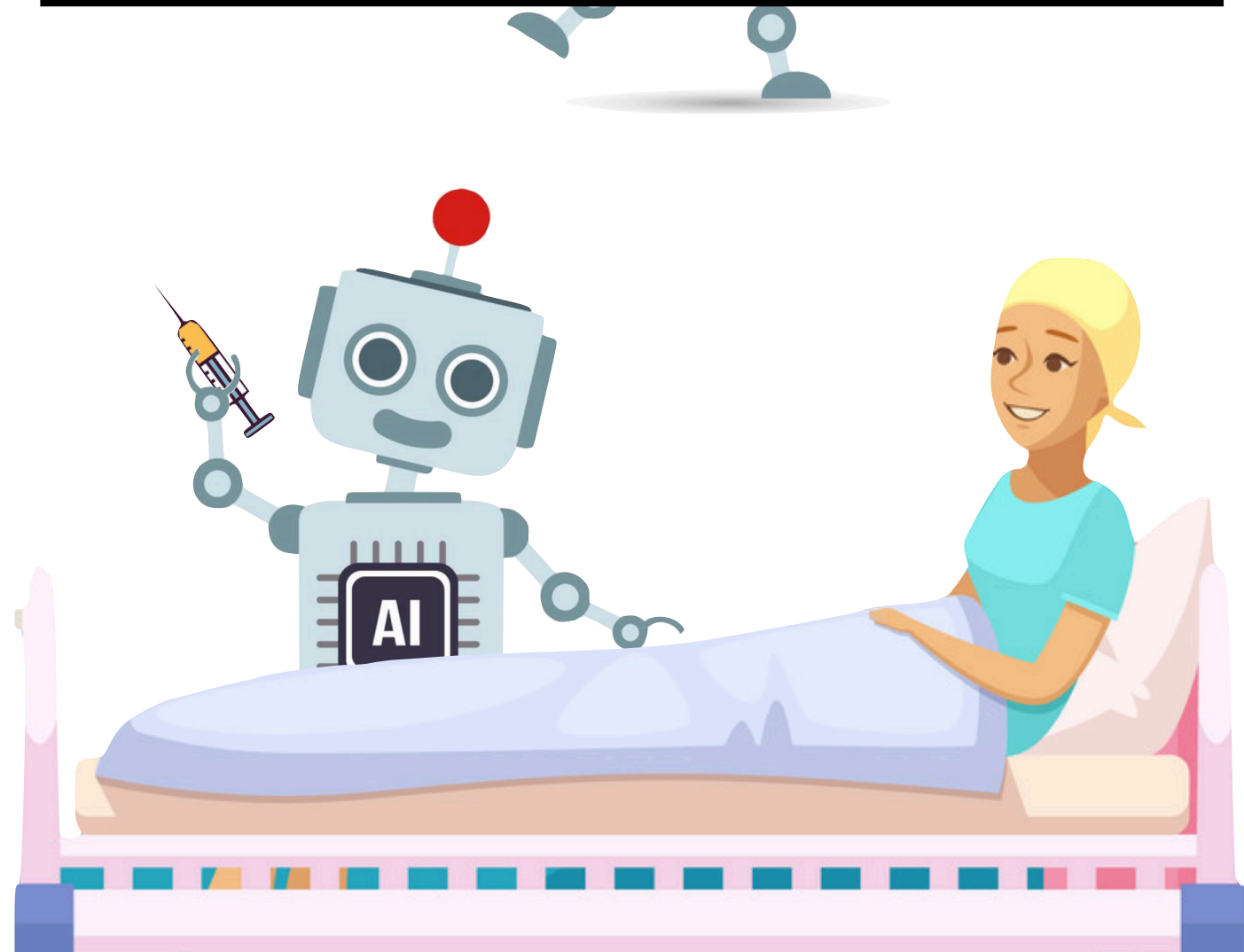
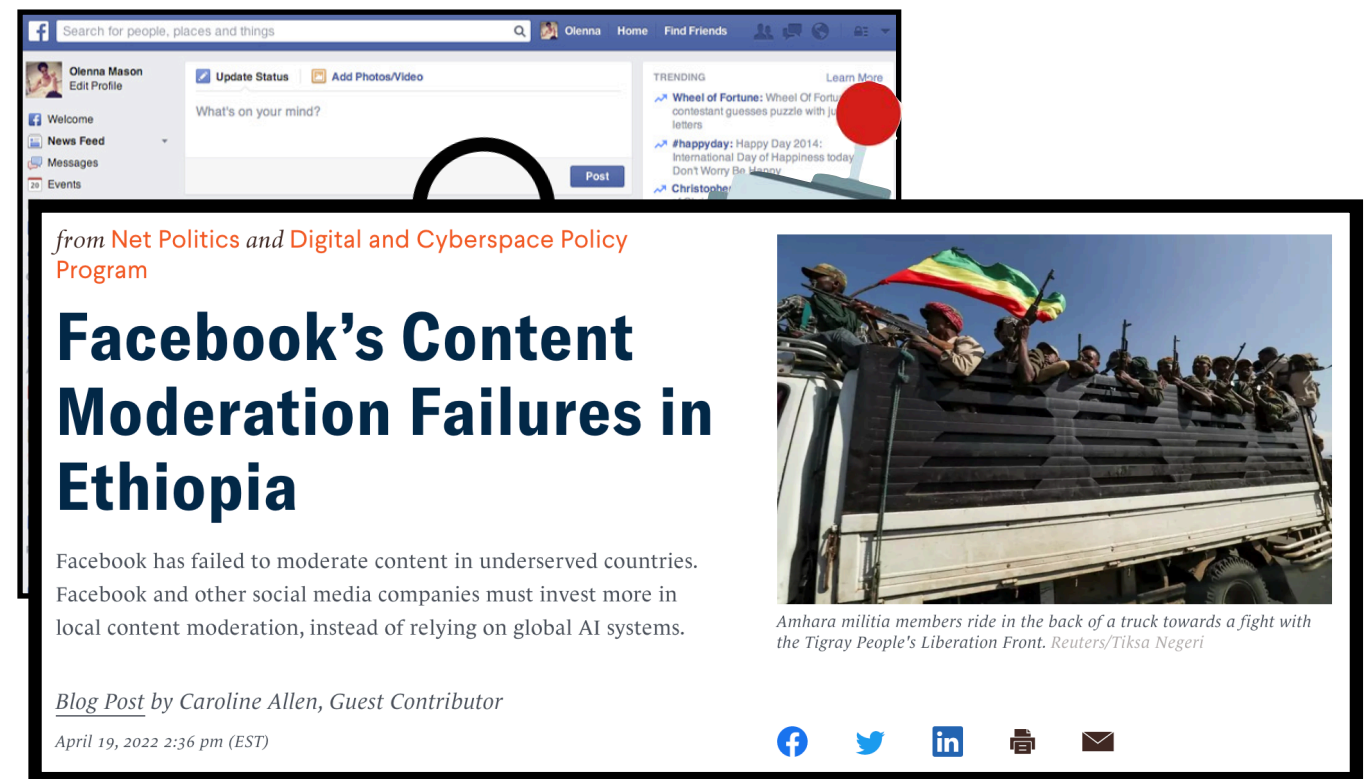
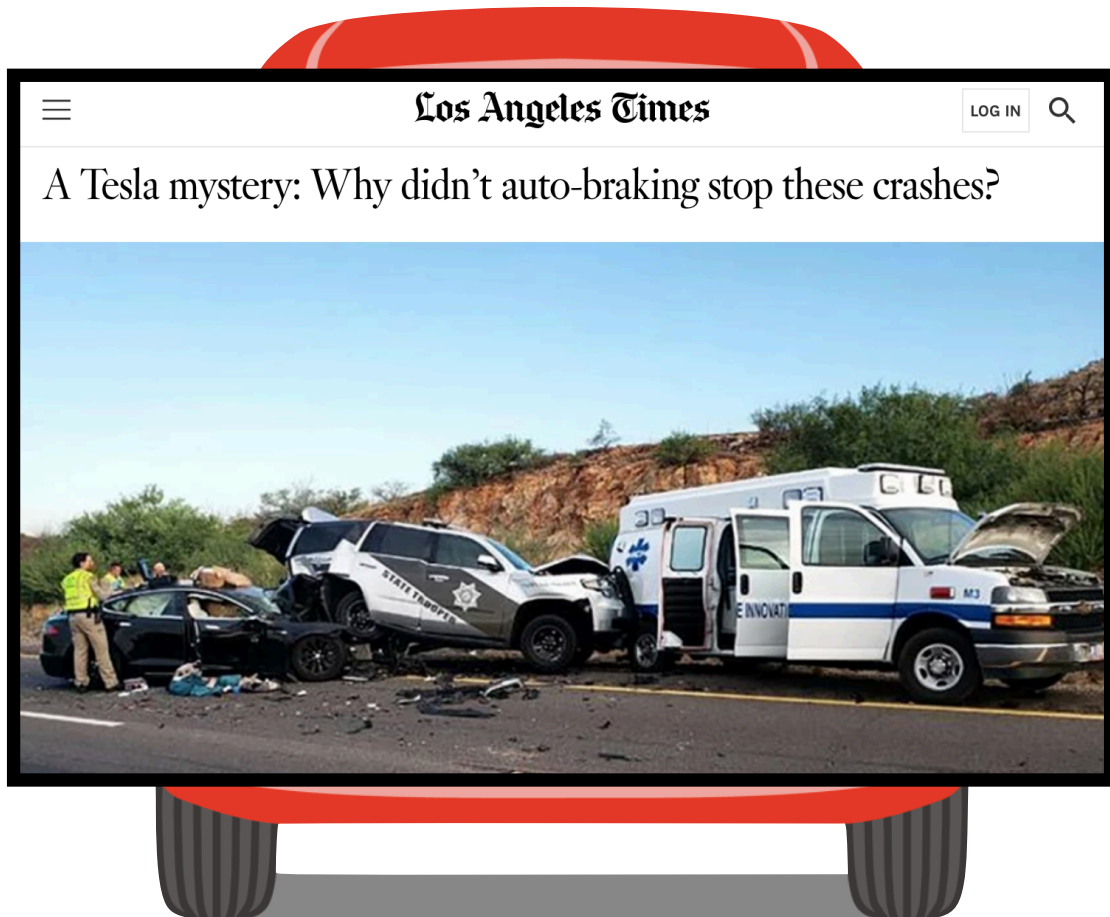


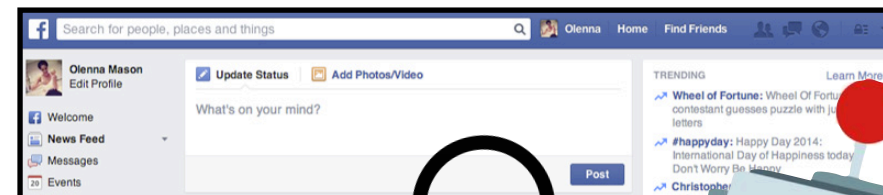












from Net Politics and Digital and Cyberspace Policy Program

## Facebook's Content Moderation Failures in Ethiopia

Facebook has failed to moderate content in underserved countries. Facebook and other social media companies must invest more in local content moderation, instead of relying on global AI systems.

*Blog Post by Caroline Allen, Guest Contributor*

April 19, 2022 2:36 pm (EST)



Amhara militia members ride in the back of a truck towards a fight with the Tigray People's Liberation Front. Reuters/Tiksa Negeri



### ARTIFICIAL INTELLIGENCE

## Google's medical AI was super accurate in a lab. Real life was a different story.

If AI is really going to make a difference, it needs to work when real humans are involved.

By Will Douglas Heaven

Medscape

Tuesday, December 13, 2022

NEWS & PERSPECTIVE

DRUGS & DISEASES

CME & EDUCATION

ACADEMY

VIDEO

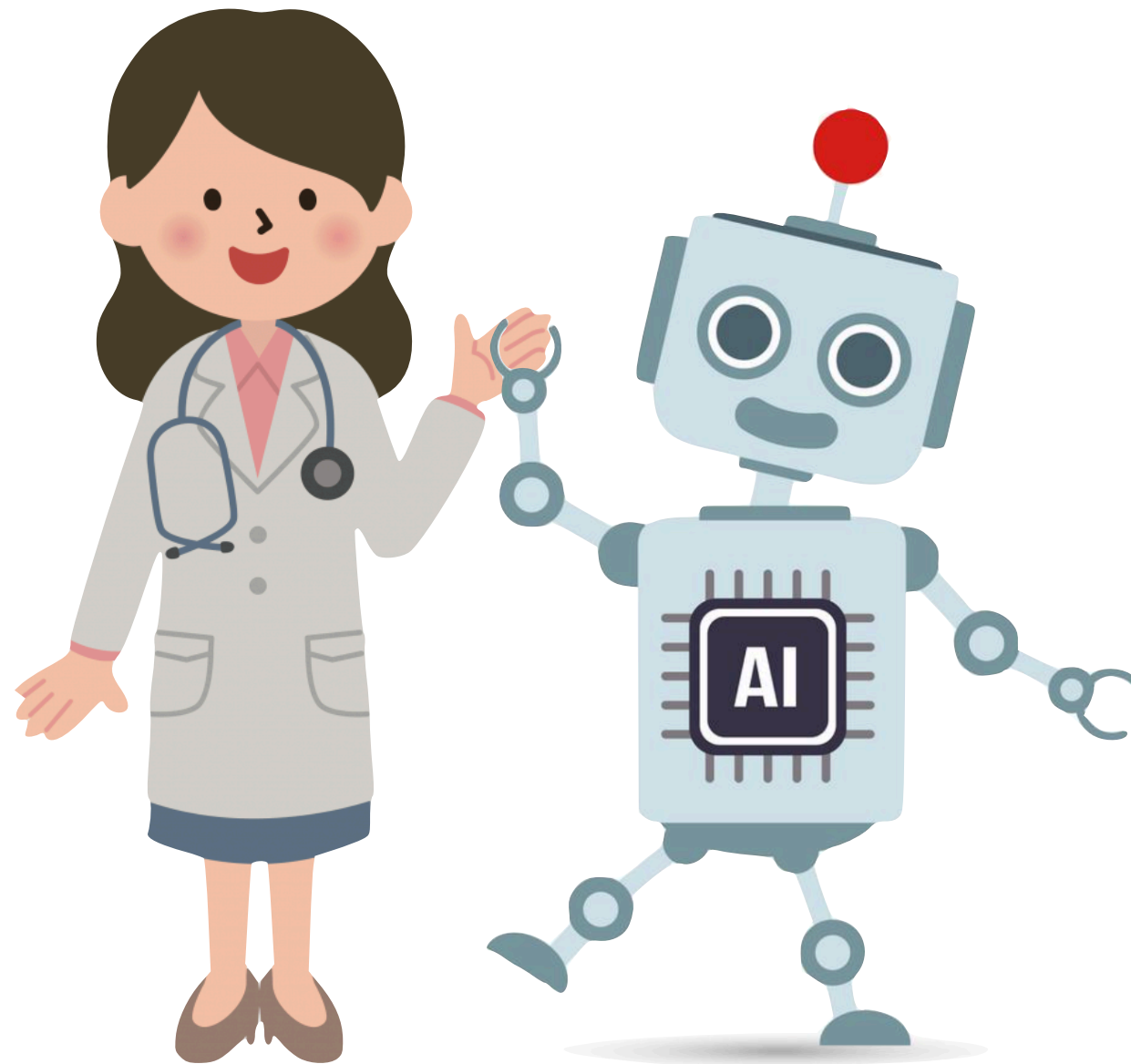
DECISION POINT

News > Medscape Medical News > Conference News > CHEST 2022

## Sepsis Predictor Tool Falls Short in Emergency Setting

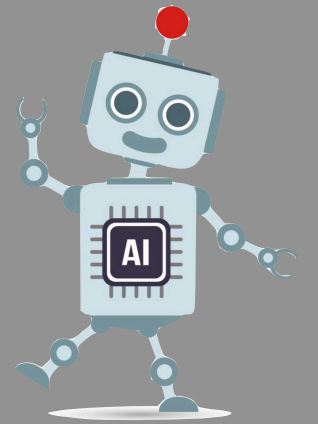
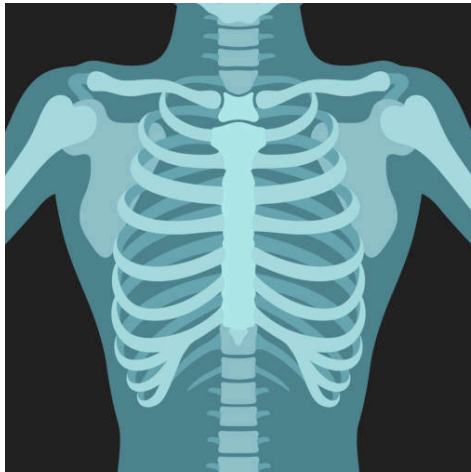
Heidi Splete

October 17, 2022



human-AI collaboration

input  
features



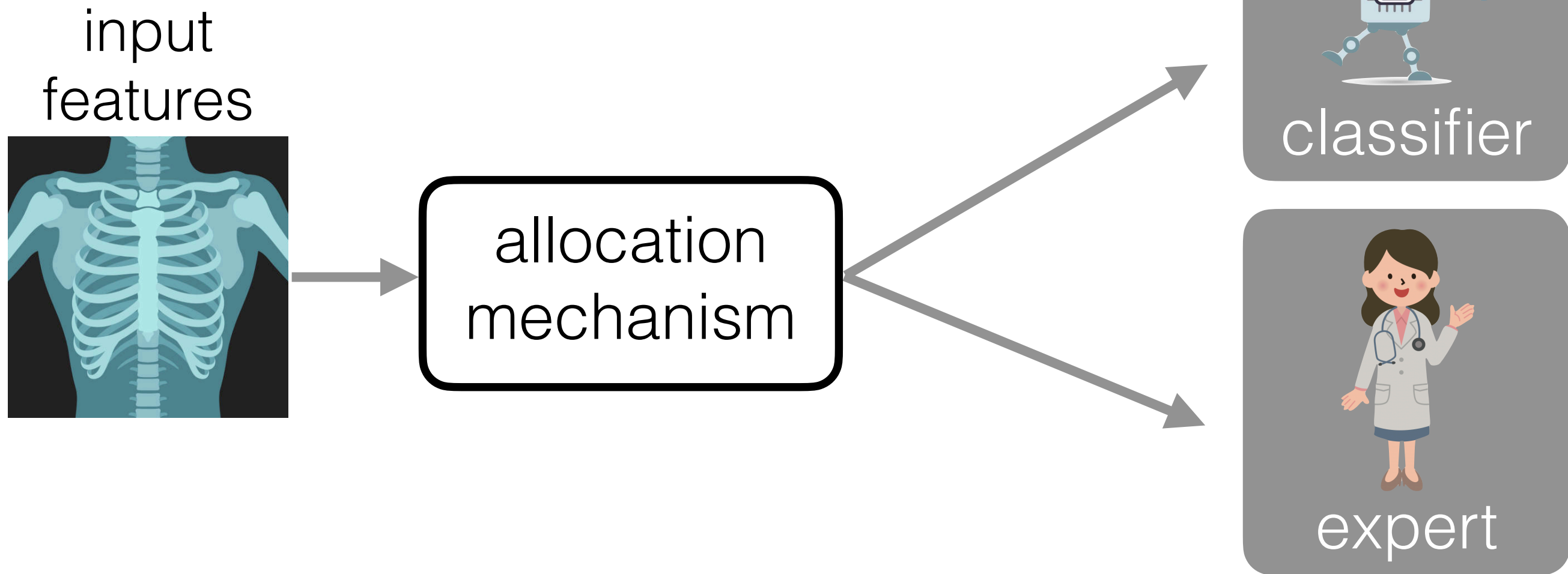
classifier



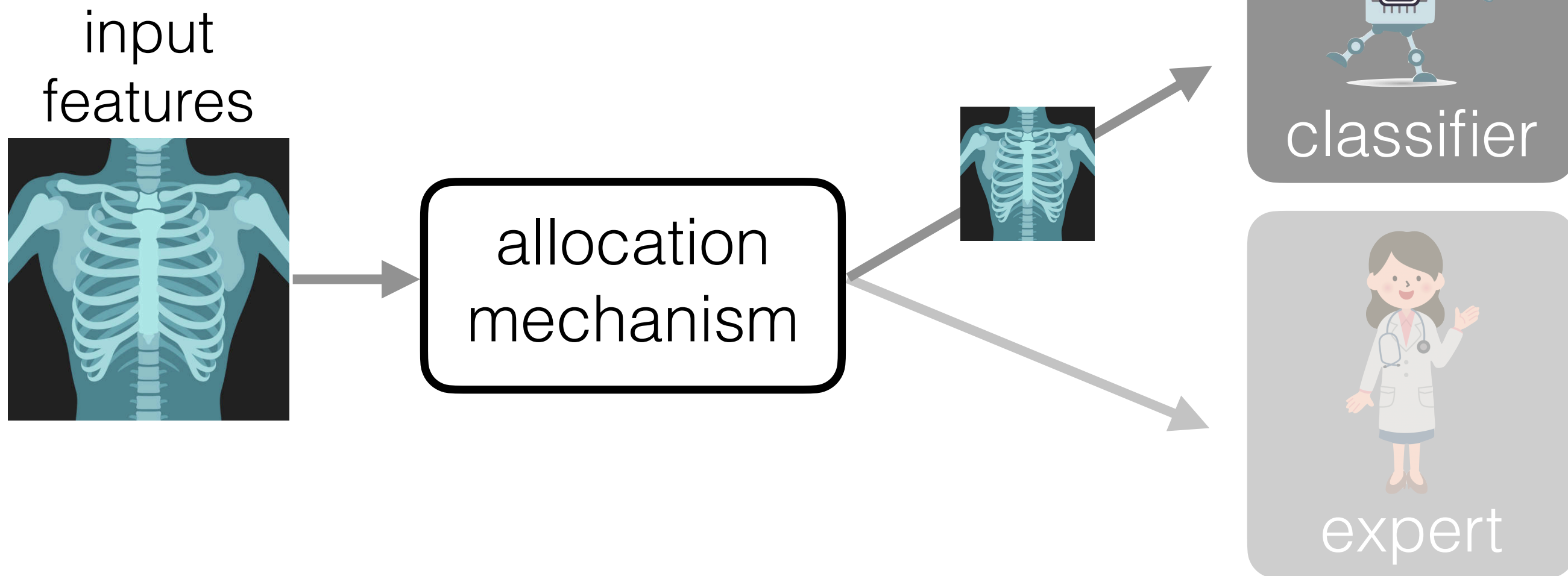
expert

learning to defer (to an expert)



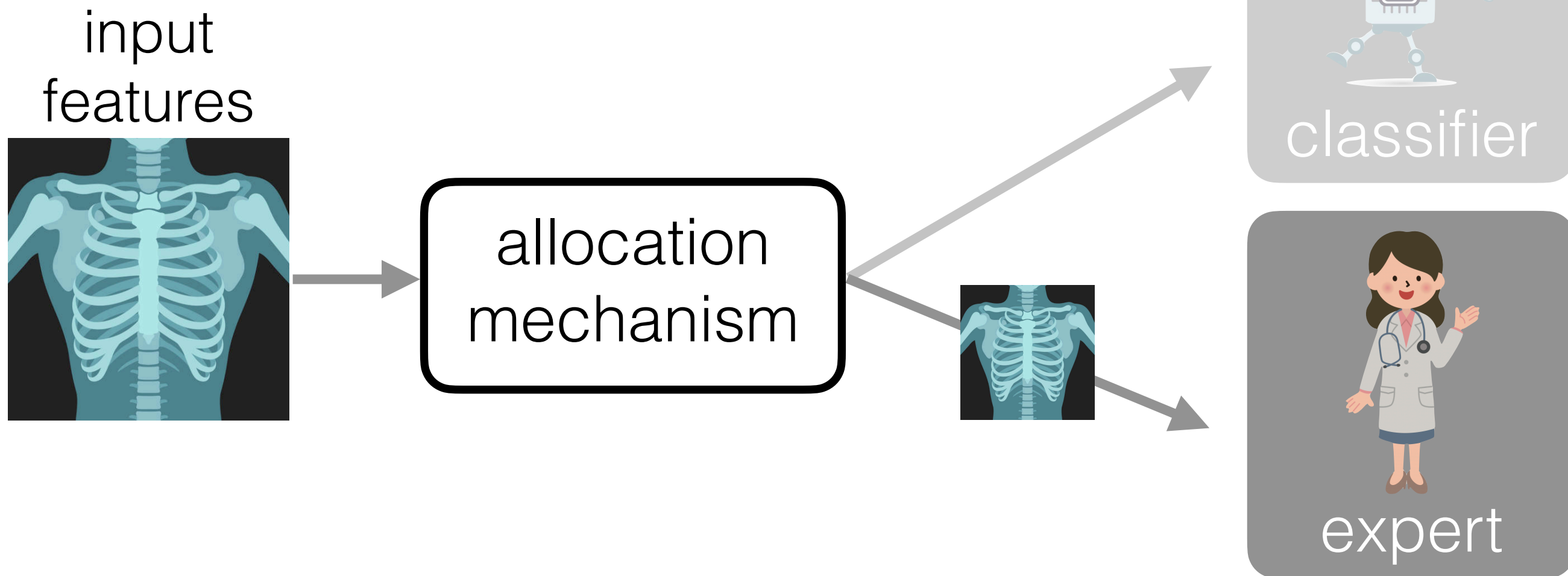


learning to defer (to an expert)

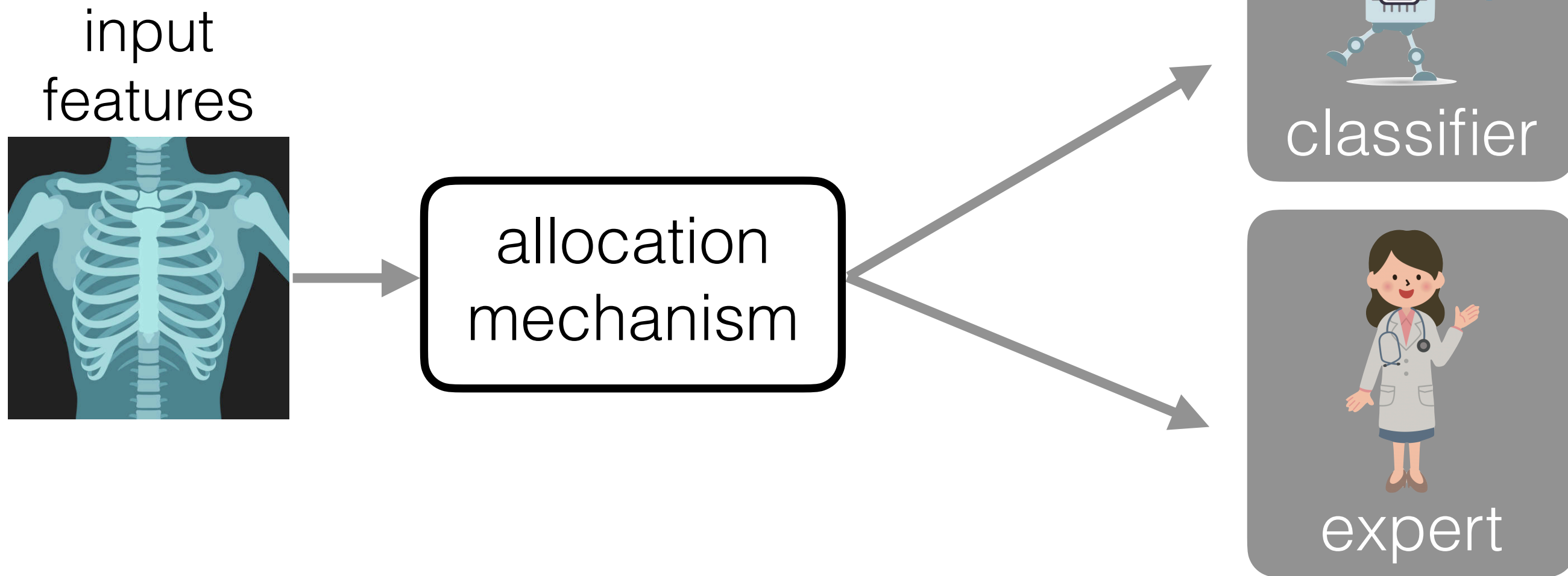


learning to defer (to an expert)



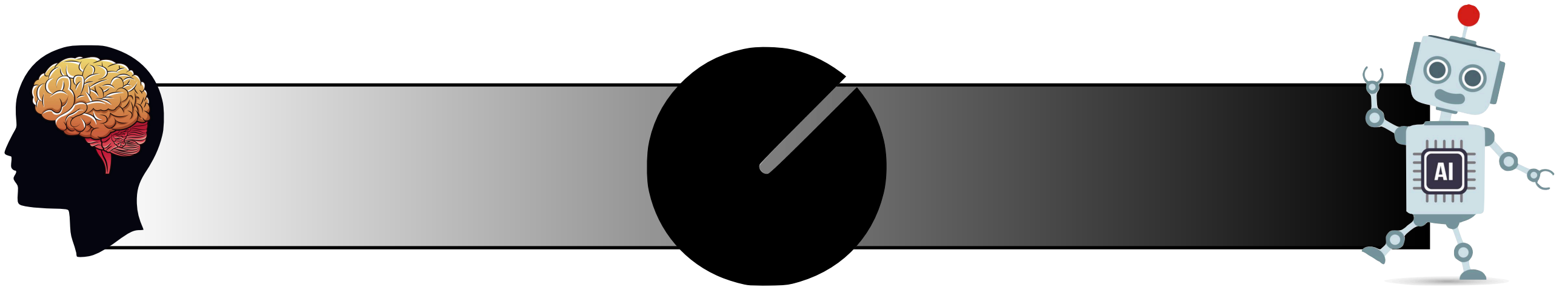


learning to defer (to an expert)



safe and robust semi-automation  
via expert handling the hardest cases

# safe, gradual automation



- ⊗ single expert
  - ⊗ softmax surrogate loss
  - ⊗ improving calibration via one-vs-all
- ⊗ multiple experts
  - ⊗ surrogate losses
  - ⊗ conformal sets of experts
- ⊗ population of experts
  - ⊗ surrogate losses
  - ⊗ meta-learning a rejector

- ⊗ **single expert**
  - ⊗ softmax surrogate loss
  - ⊗ improving calibration via one-vs-all
  
- ⊗ **multiple experts**
  - ⊗ surrogate losses
  - ⊗ conformal sets of experts
  
- ⊗ **population of experts**
  - ⊗ surrogate losses
  - ⊗ meta-learning a rejector

## ⊗ single expert

- ⊗ softmax surrogate loss

- ⊗ improving calibration via one-vs-all

## ⊗ multiple experts

- ⊗ surrogate losses

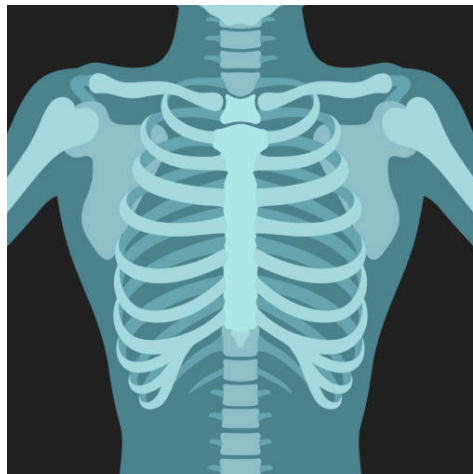
- ⊗ conformal sets of experts

## ⊗ population of experts

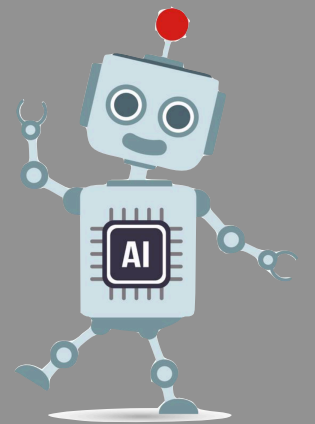
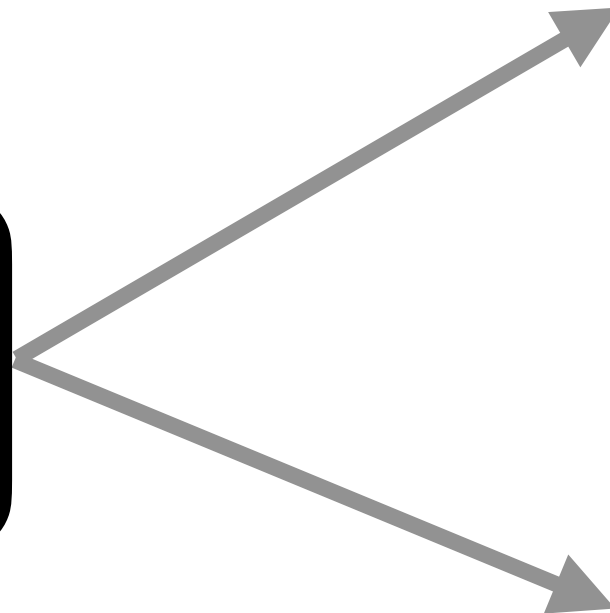
- ⊗ surrogate losses

- ⊗ meta-learning a rejector

input  
features



allocation  
mechanism

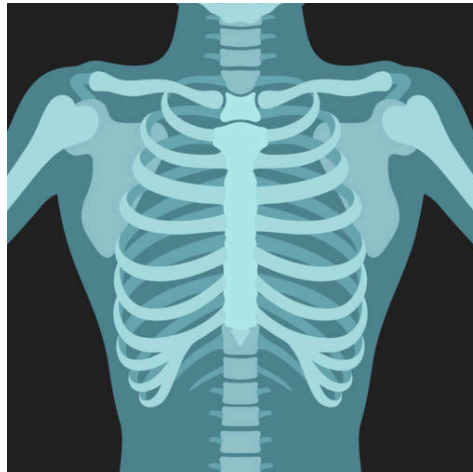


classifier



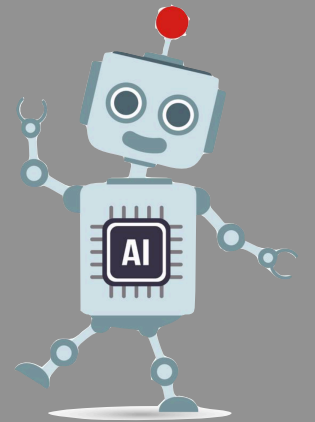
expert

input  
features



allocation  
mechanism

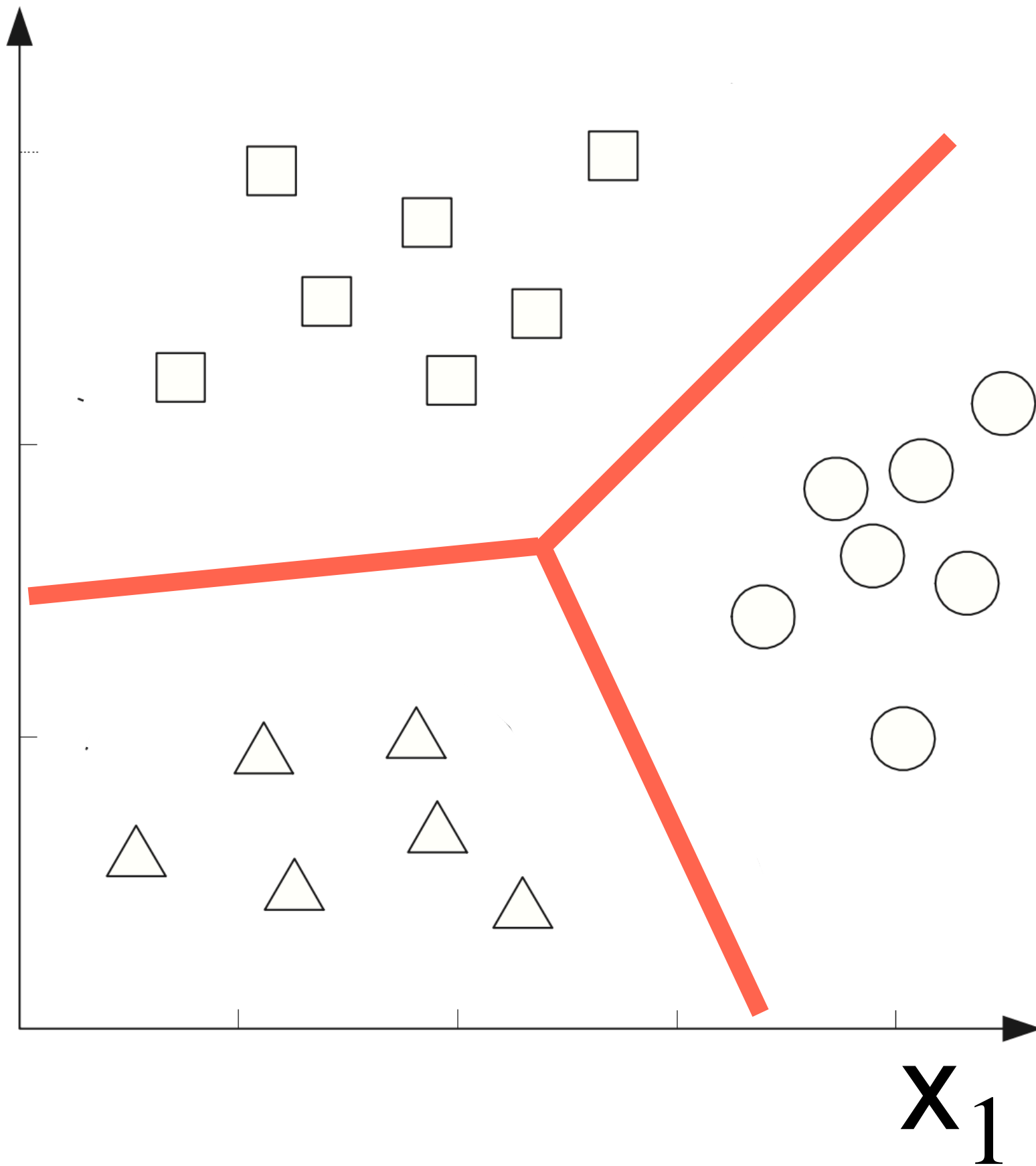
classifier



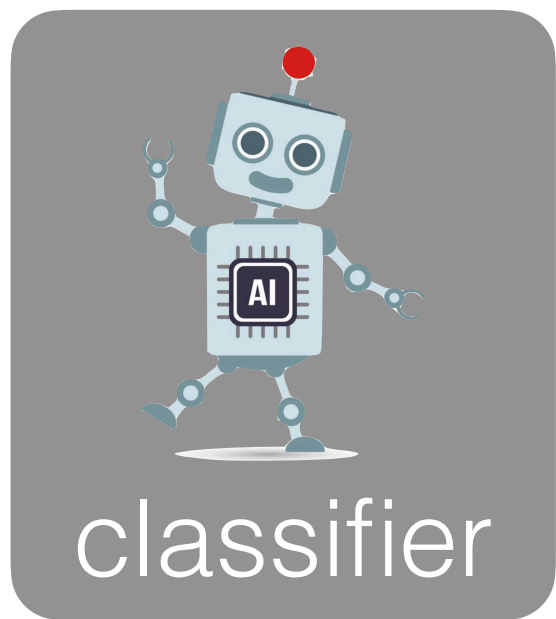
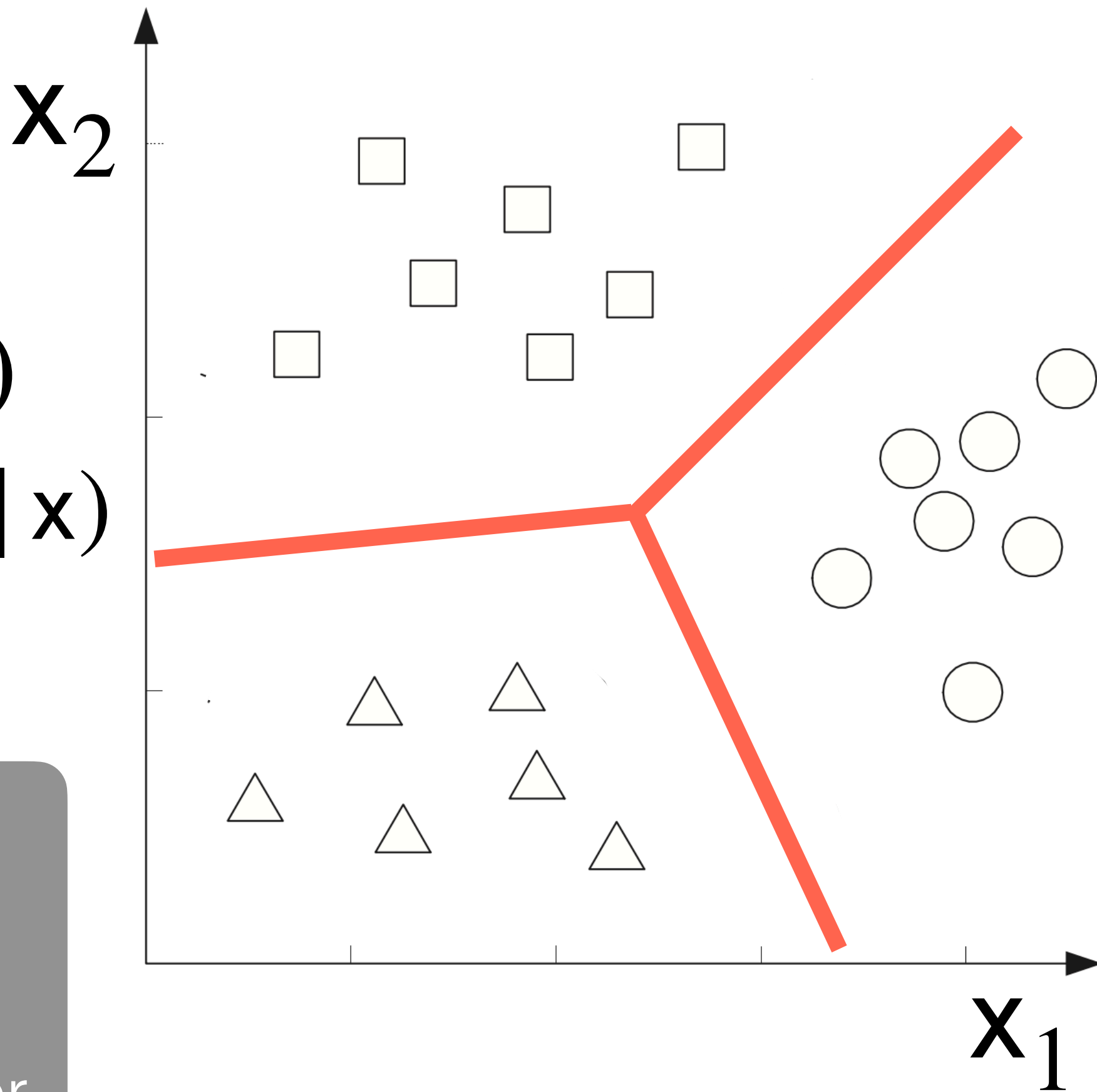
expert



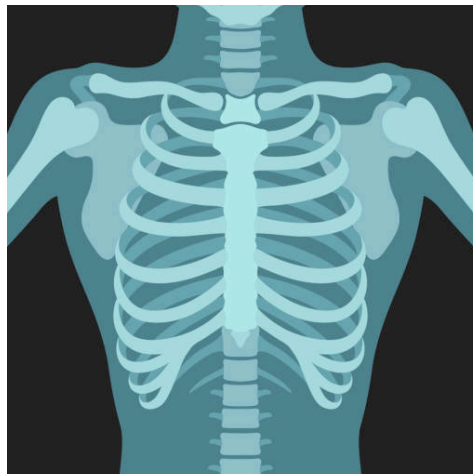


$$x_2$$


$$p(y | x) \approx \mathbb{P}(y | x)$$

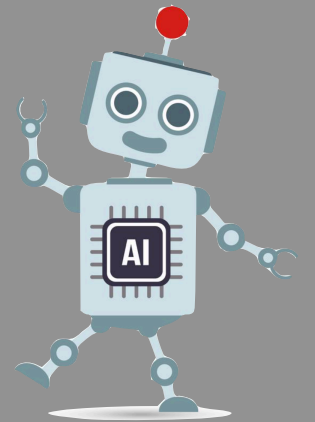


input  
features



allocation  
mechanism

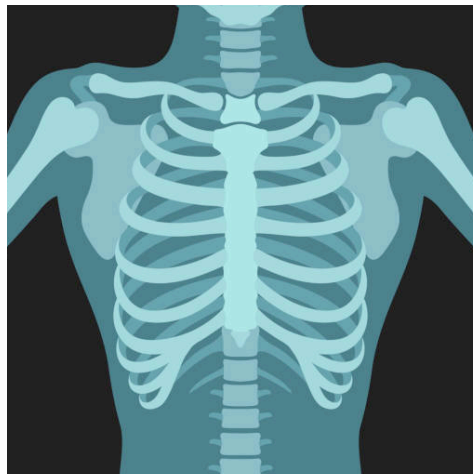
classifier



expert

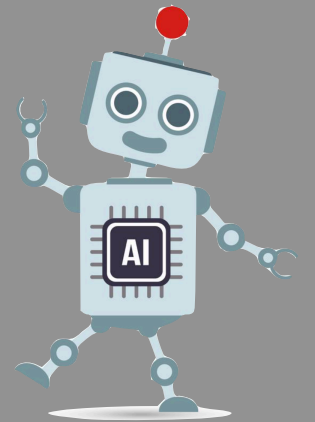


input  
features



allocation  
mechanism

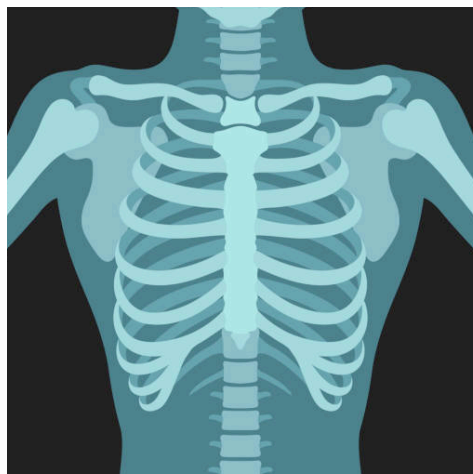
classifier



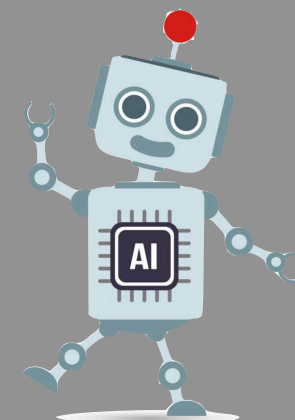
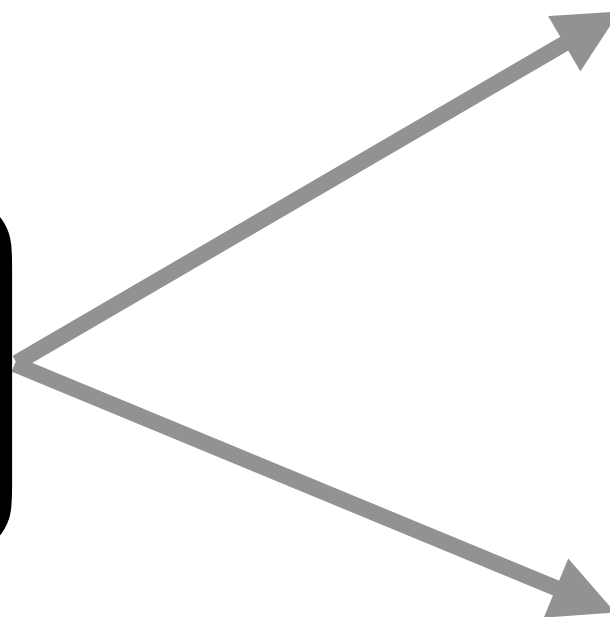
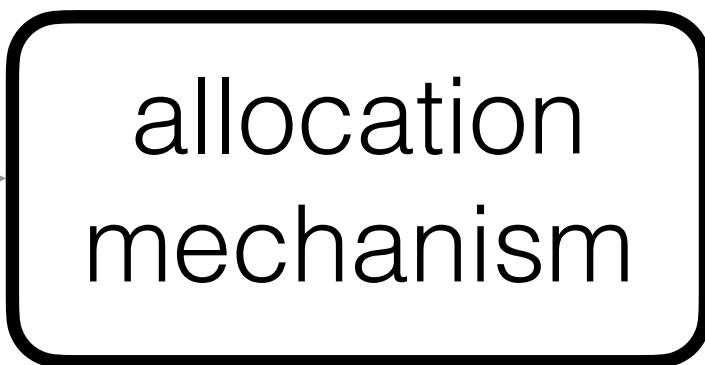
expert



input  
features



allocation  
mechanism



classifier



expert

input  
features

**X**

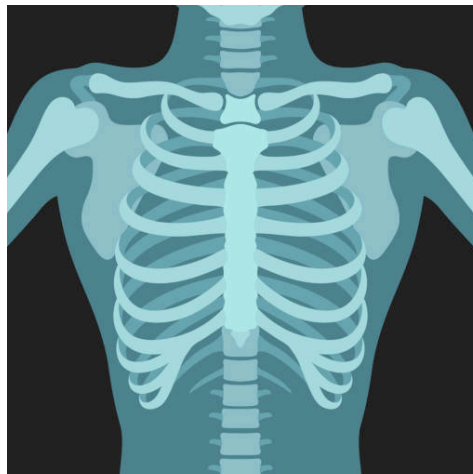


(black box)

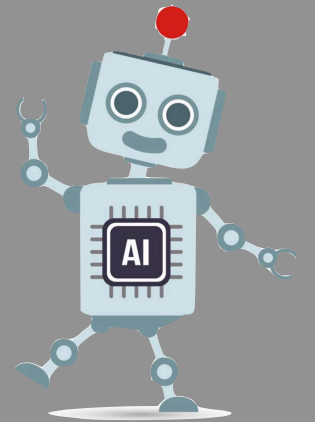
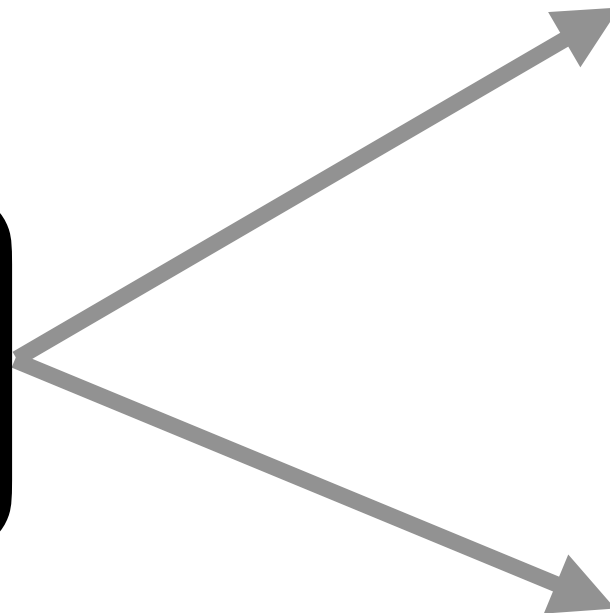


**m**  
prediction

input  
features



allocation  
mechanism

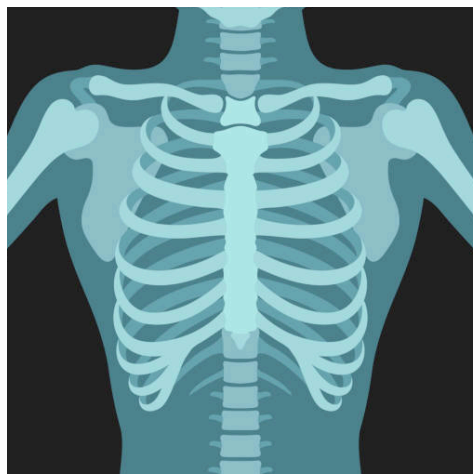


classifier



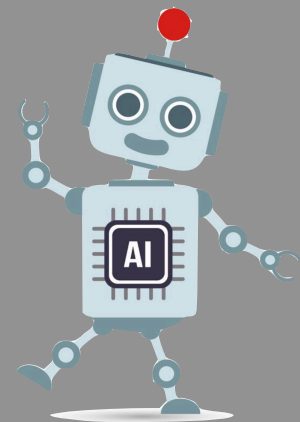
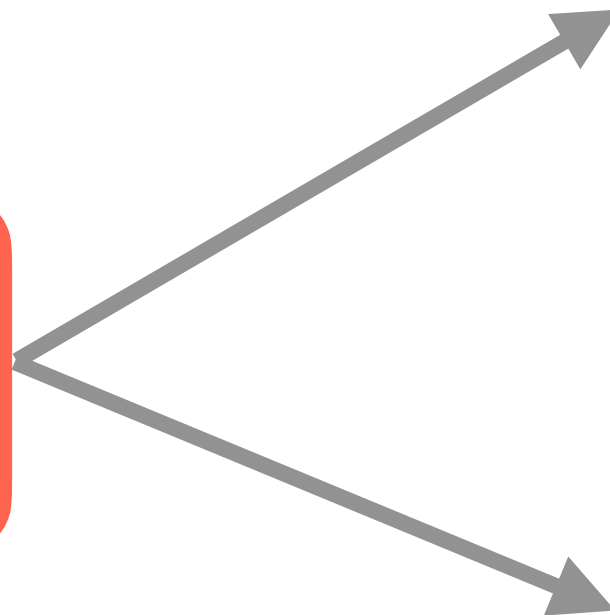
expert

input  
features



allocation  
mechanism

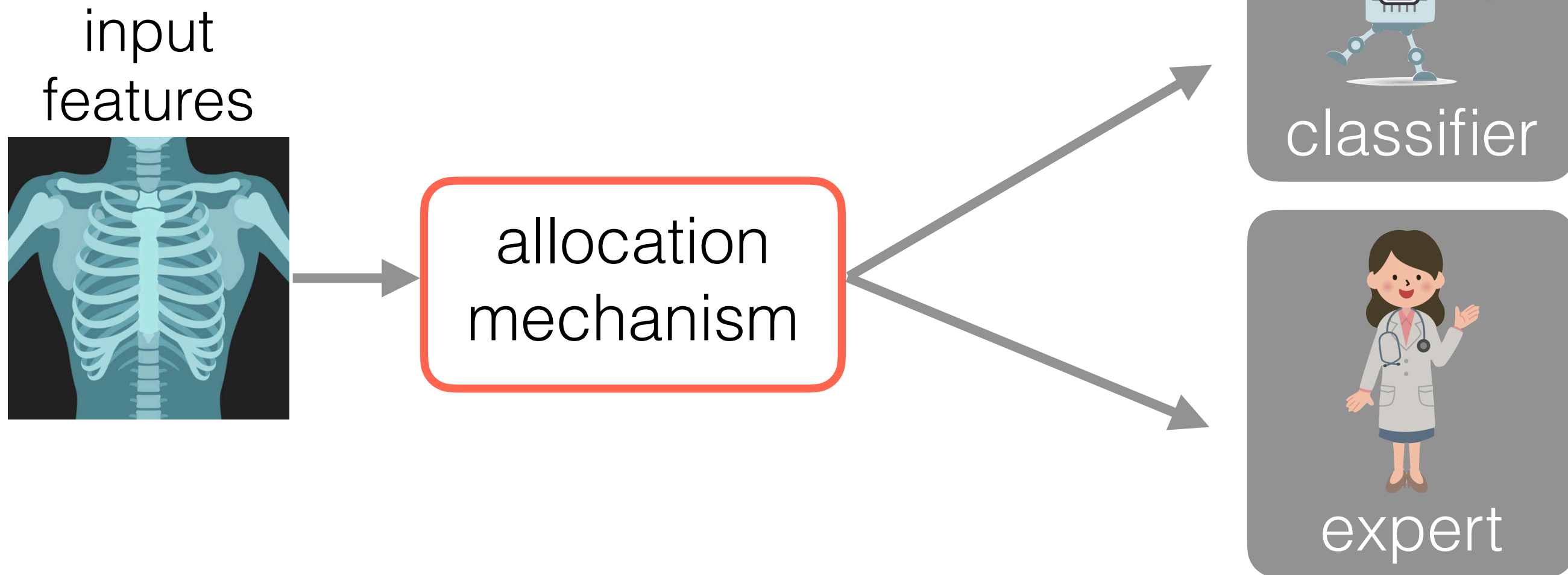
???



classifier



expert

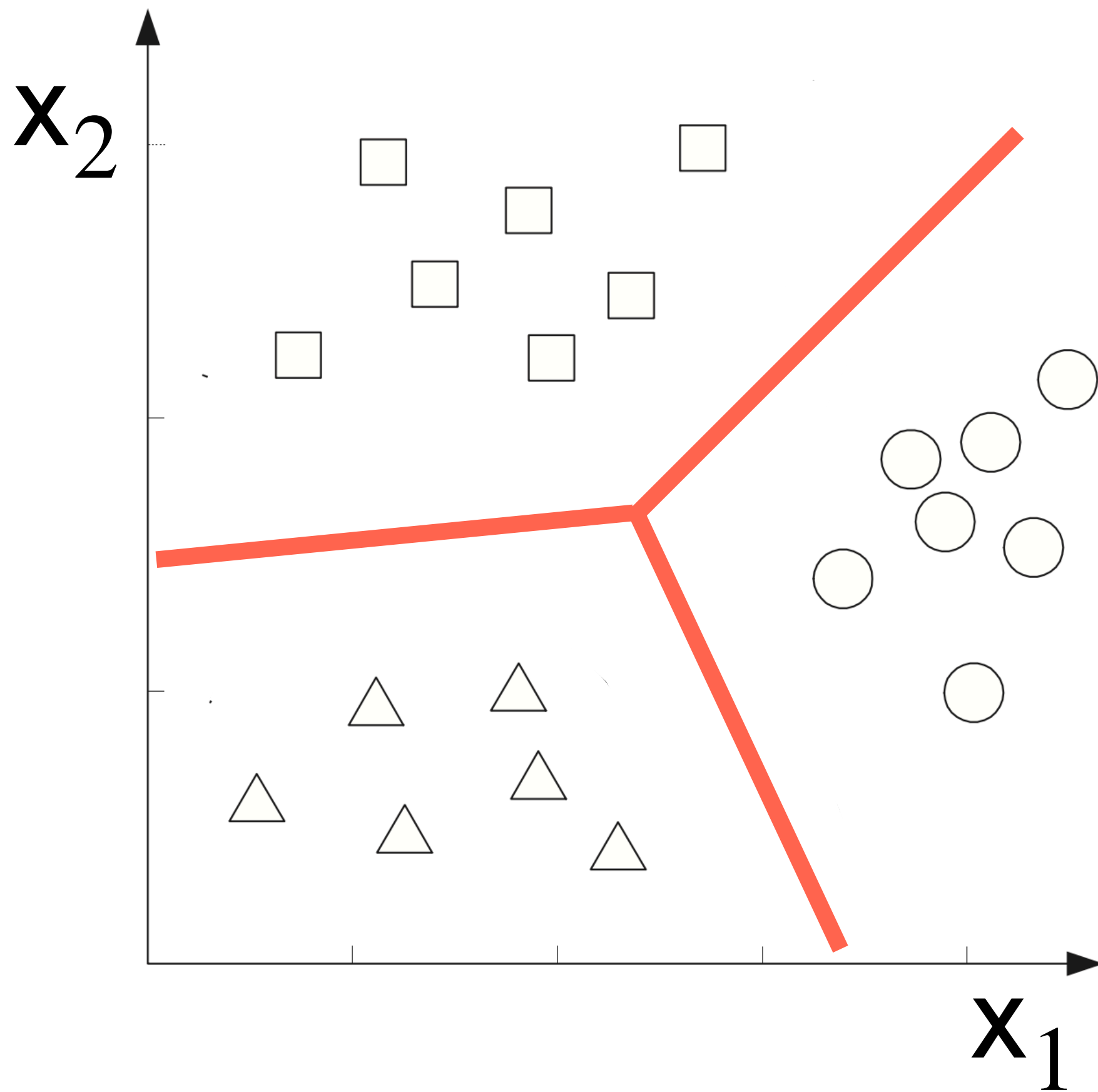


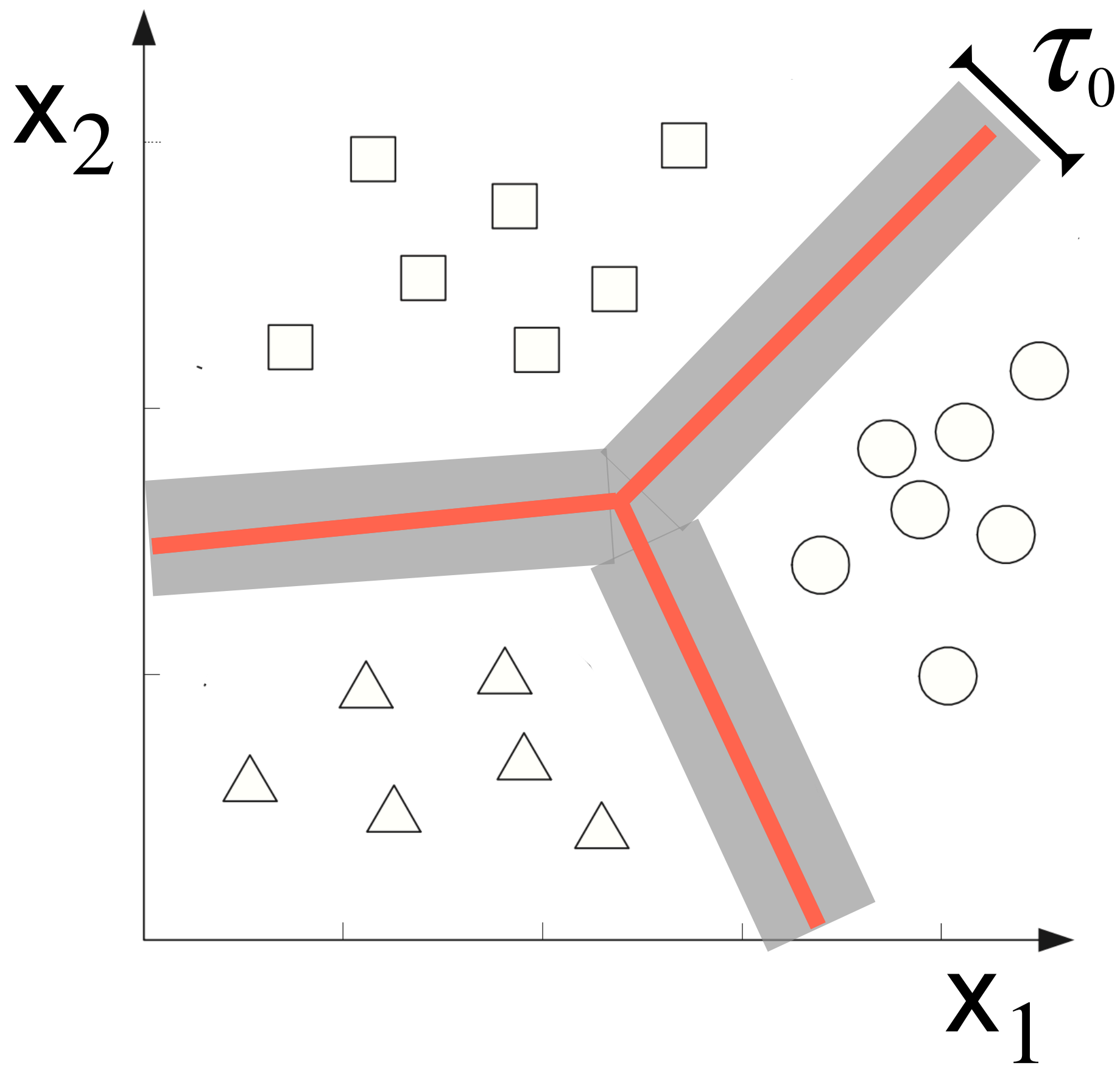
defer to expert if...

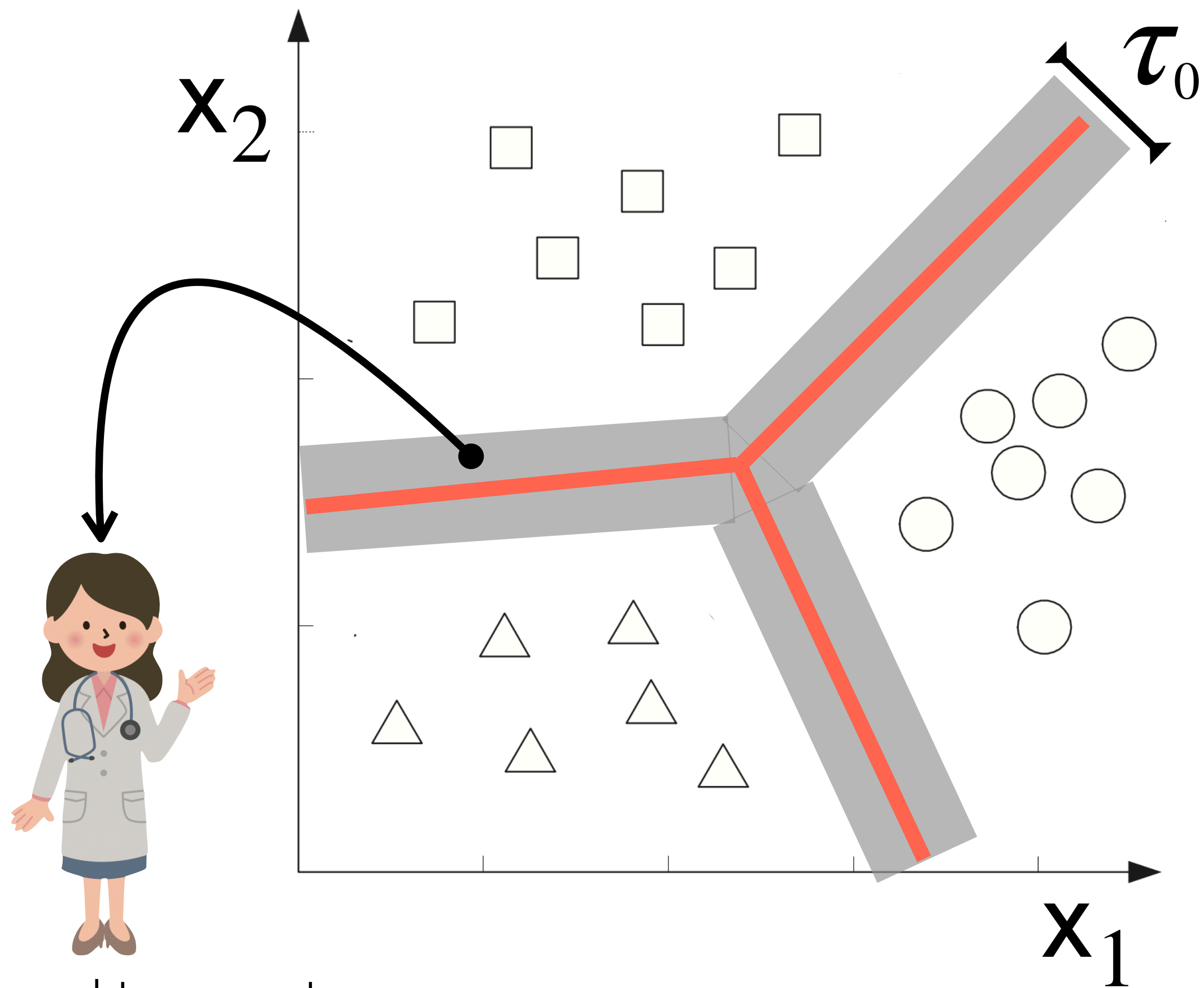
$$\max_y p(y | \mathbf{x}) \leq \tau_0$$

(constant)

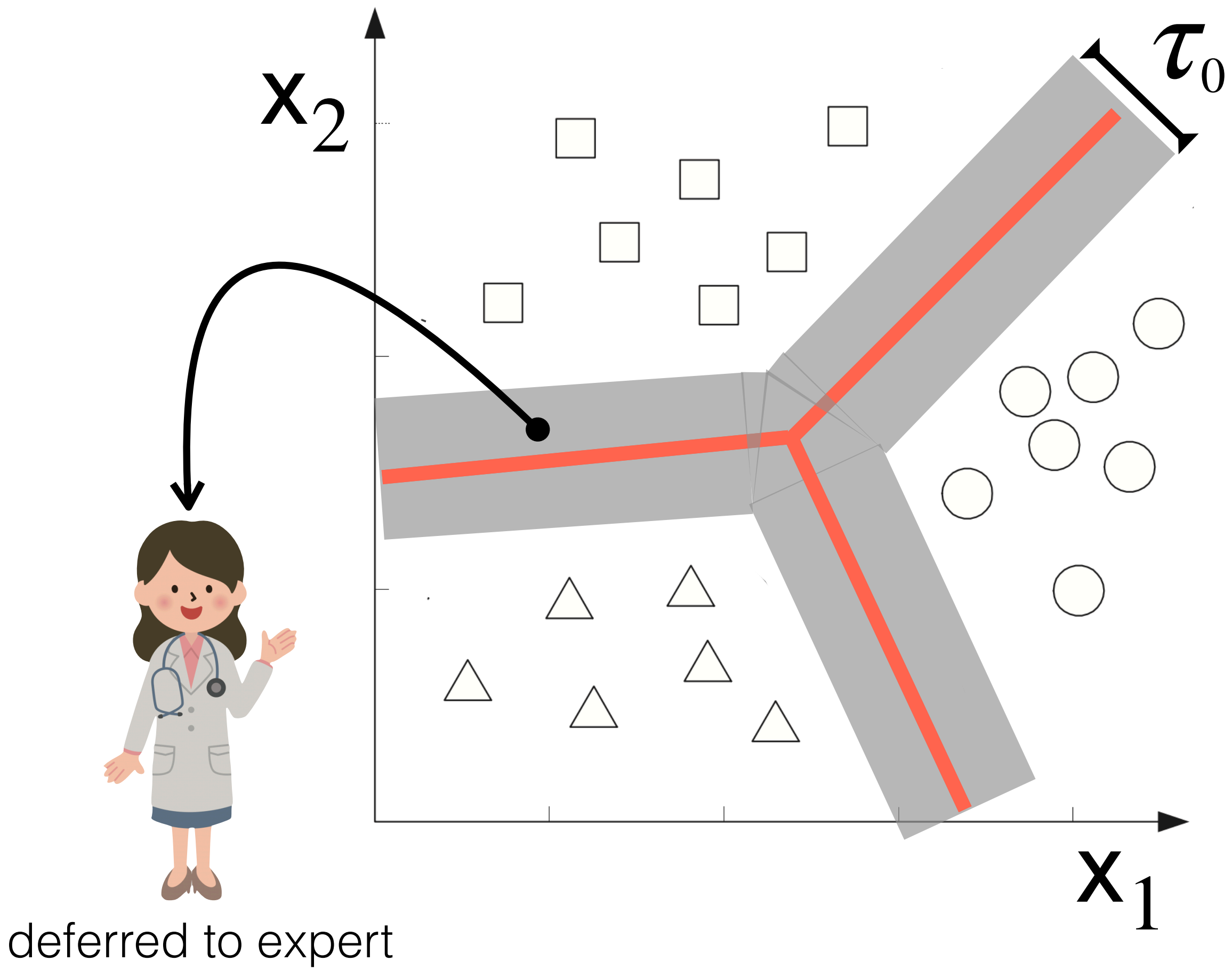


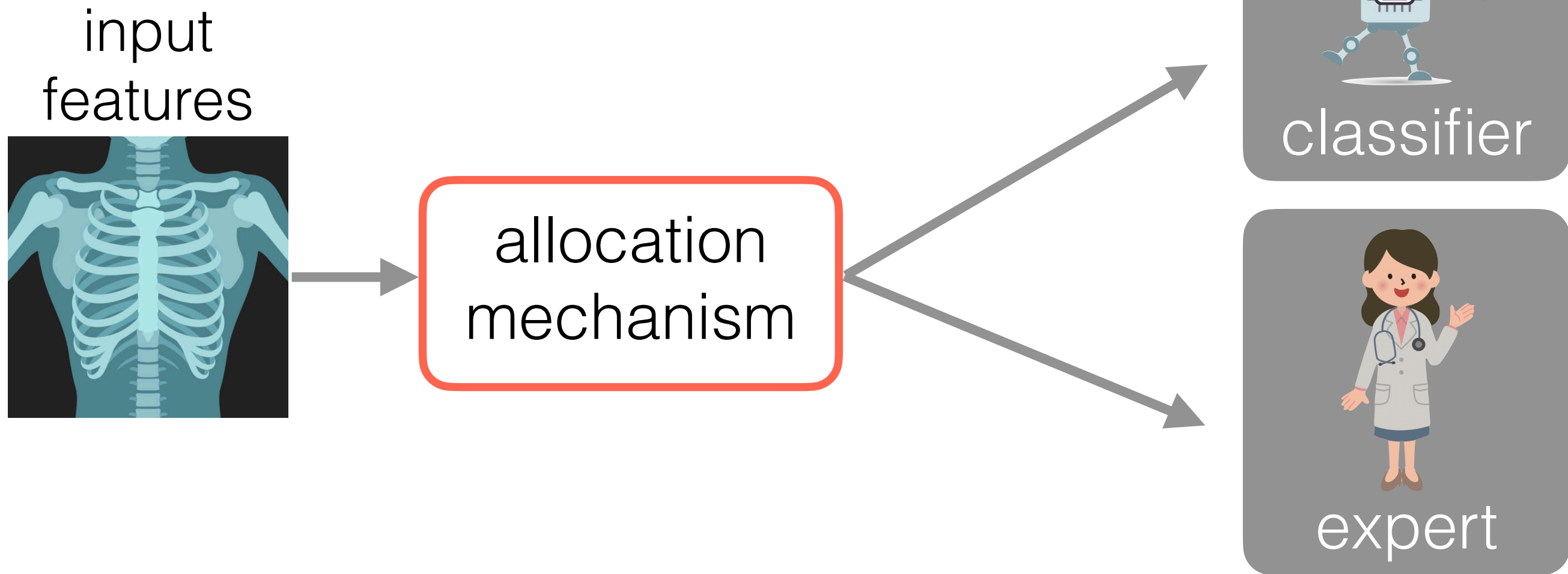






deferred to expert

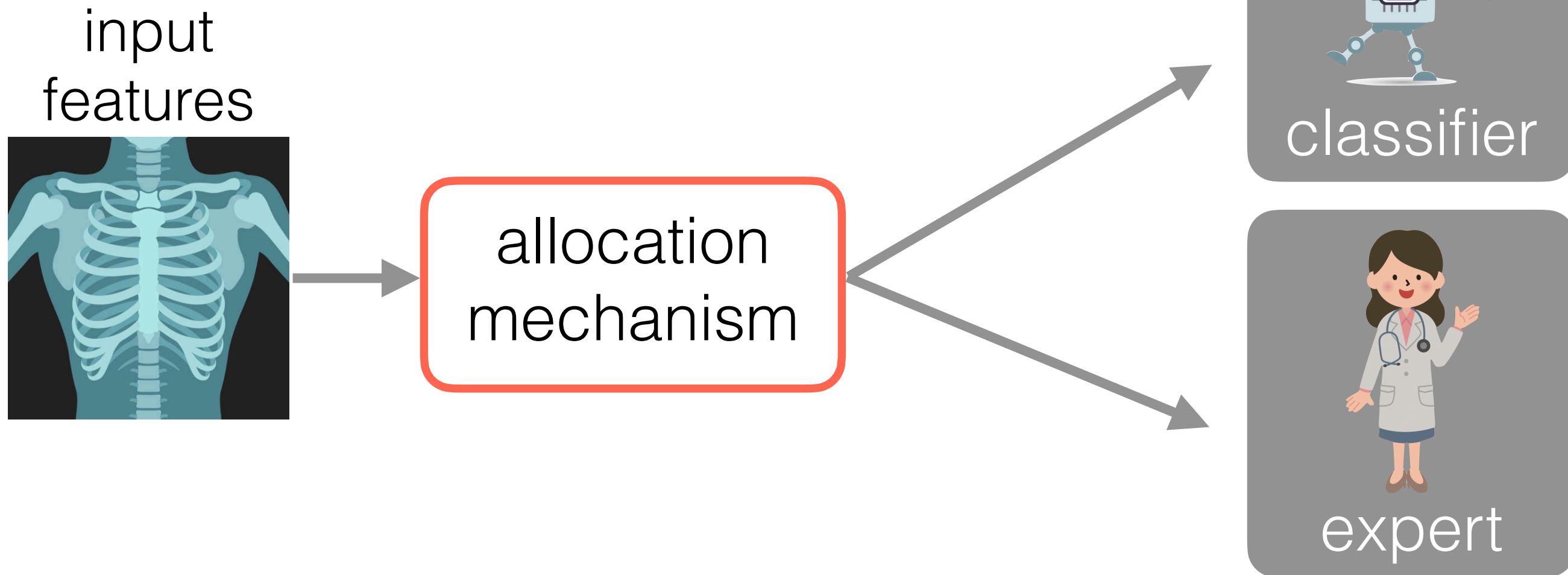




defer to expert if...

$$\max_y p(y | \mathbf{x}) \leq \tau_0$$

(constant)

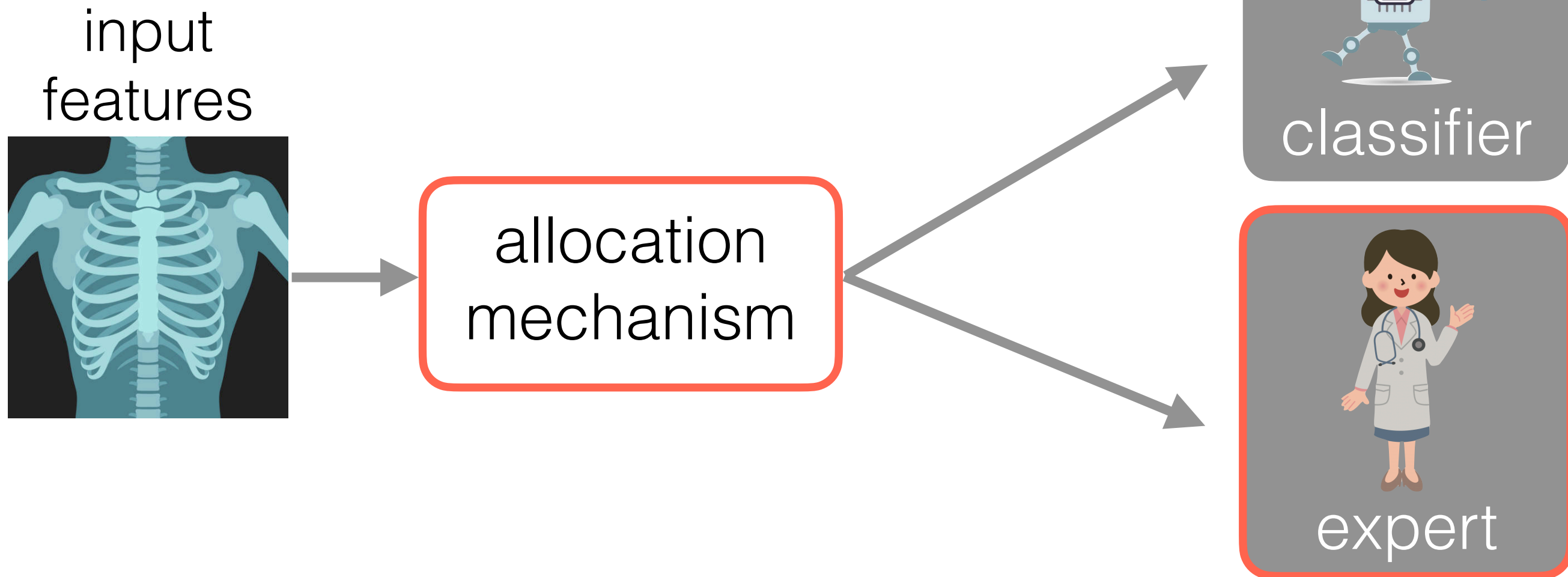


defer to expert if...

$$\max_y p(y | \mathbf{x}) \leq \tau_0$$

(constant)

problem?



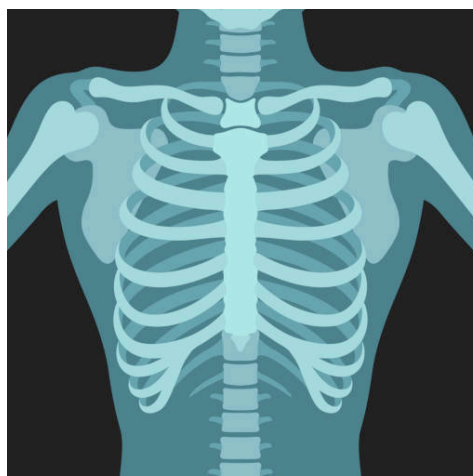
defer to expert if...

$$\max_y p(y | \mathbf{x}) \leq \tau_0$$

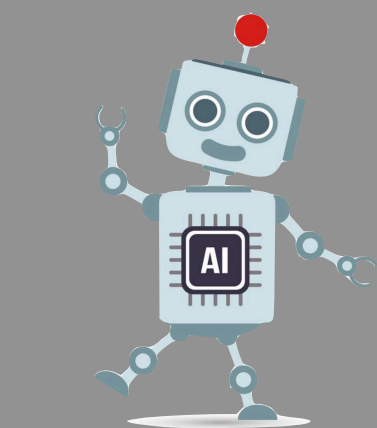
(constant)

the expert's  
knowledge is  
not considered!

input  
features



allocation  
mechanism



classifier

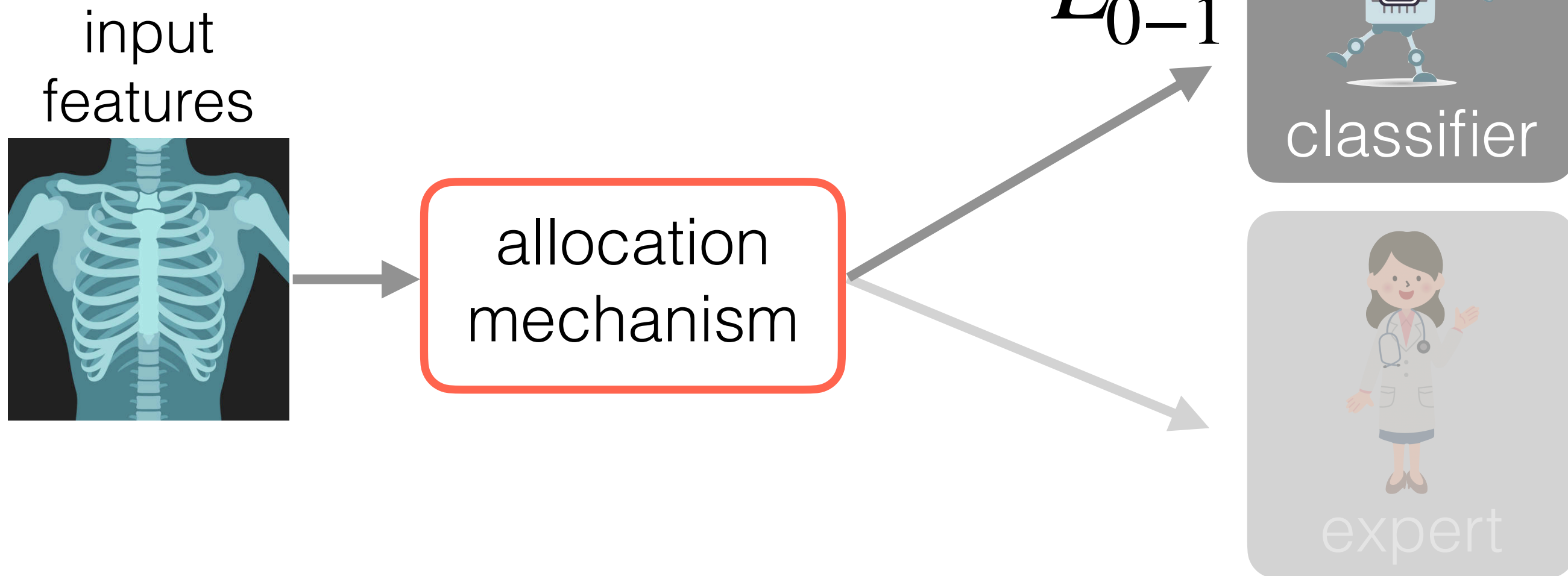


expert

defer to expert if...

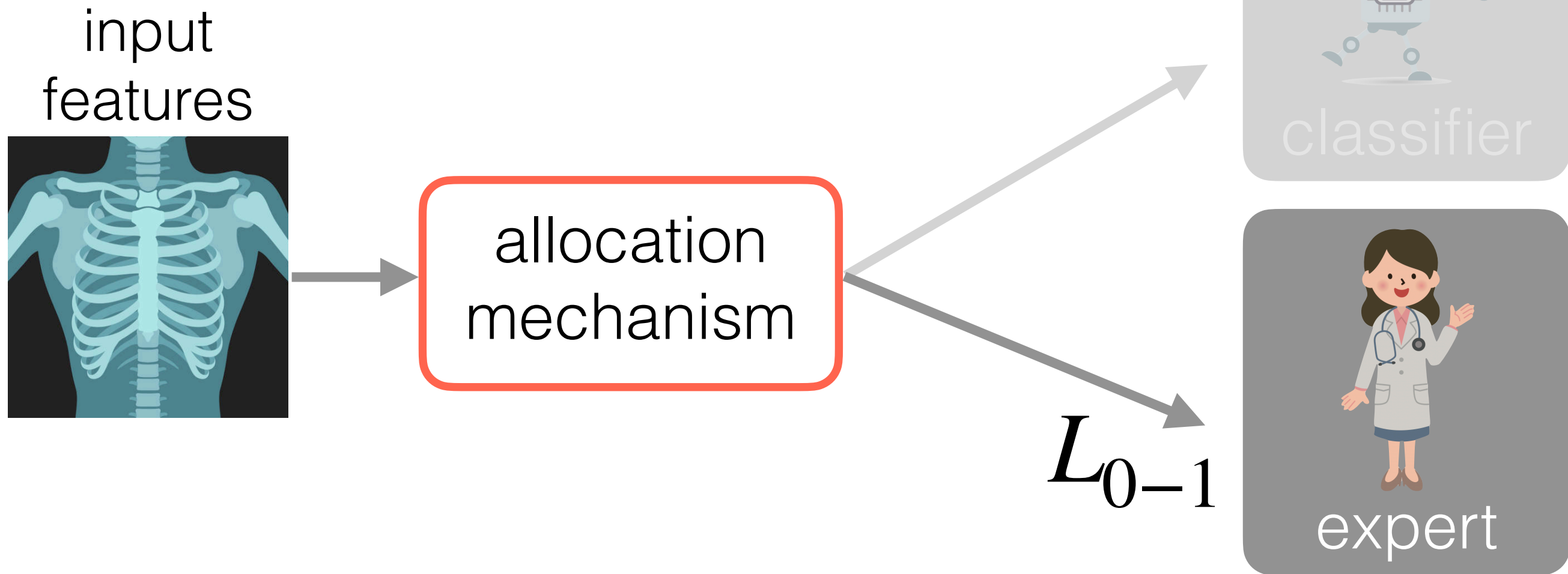
$$\max_y p(y | \mathbf{x}) \leq \tau \left( \text{X-ray image}, \text{expert} \right)$$





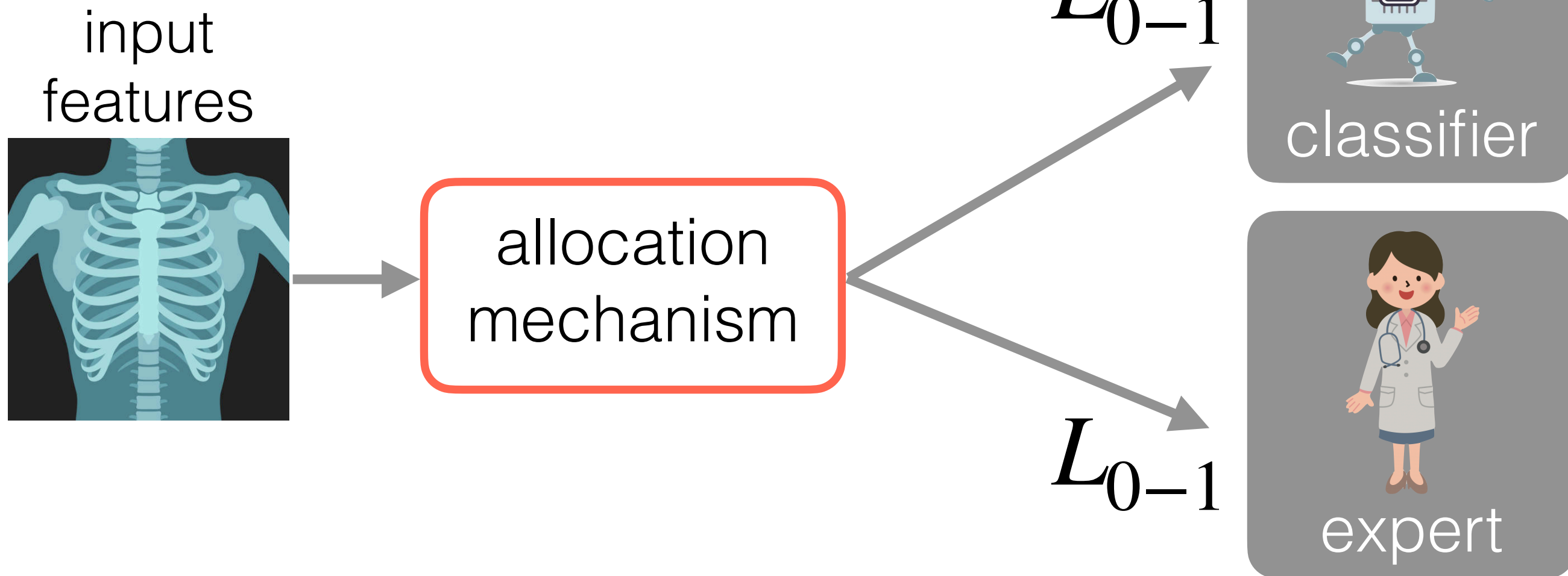
defer to expert if...

$$\max_y p(y | \mathbf{x}) \leq \tau \left( \text{input features}, \text{expert} \right)$$



defer to expert if...

$$\max_y p(y | \mathbf{x}) \leq \tau \left( \text{input features}, \text{expert} \right)$$



Bayes optimal deferral rule:

$$\max_y \mathbb{P}(y | \mathbf{x}) \leq \mathbb{P}(\mathbf{m} = y | \mathbf{x})$$

*probability that the expert is correct*

# softmax implementation

[Mozannar & Sontag, 2020]

# softmax implementation

[Mozannar & Sontag, 2020]

training data

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

# softmax implementation

[Mozannar & Sontag, 2020]

training data

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

model

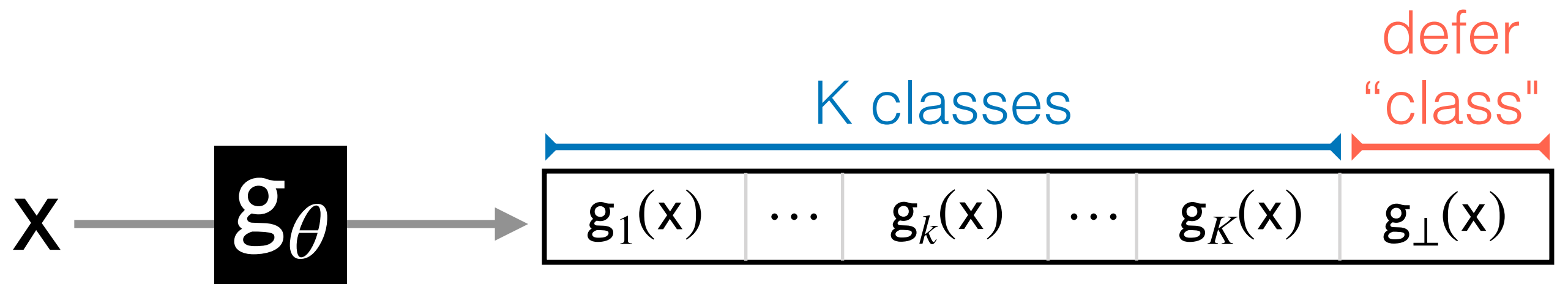
# softmax implementation

[Mozannar & Sontag, 2020]

training data

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

model



# softmax implementation

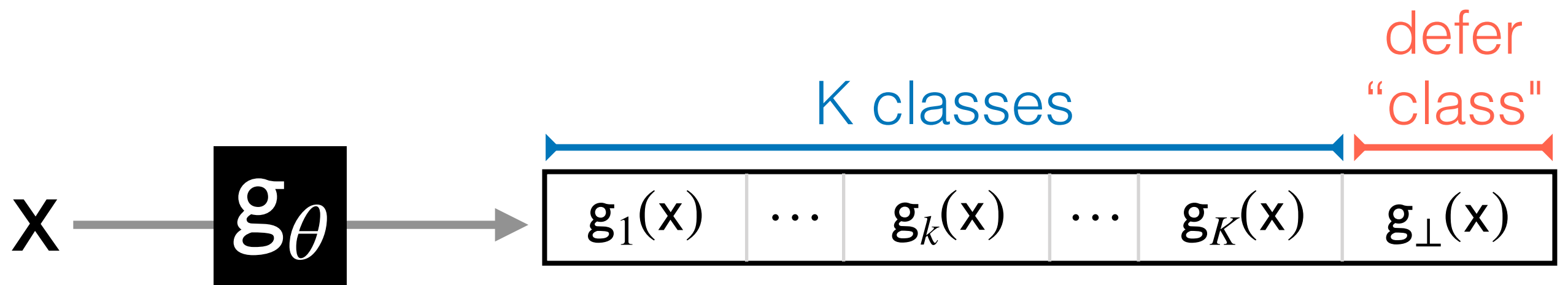
[Mozannar & Sontag, 2020]

training data

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

model

$$\mathbf{g}_k(\mathbf{x}) \in \mathbb{R}$$





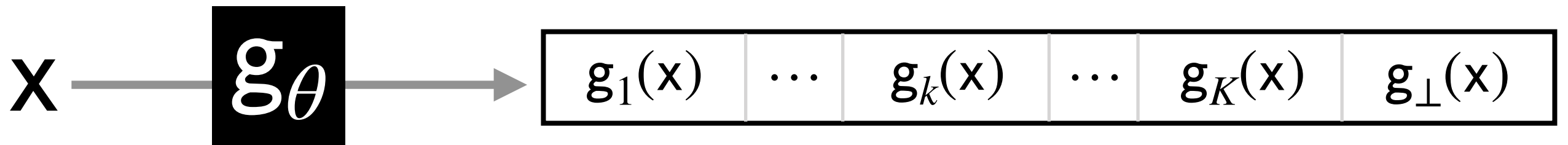
# softmax implementation

[Mozannar & Sontag, 2020]

training data

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

model



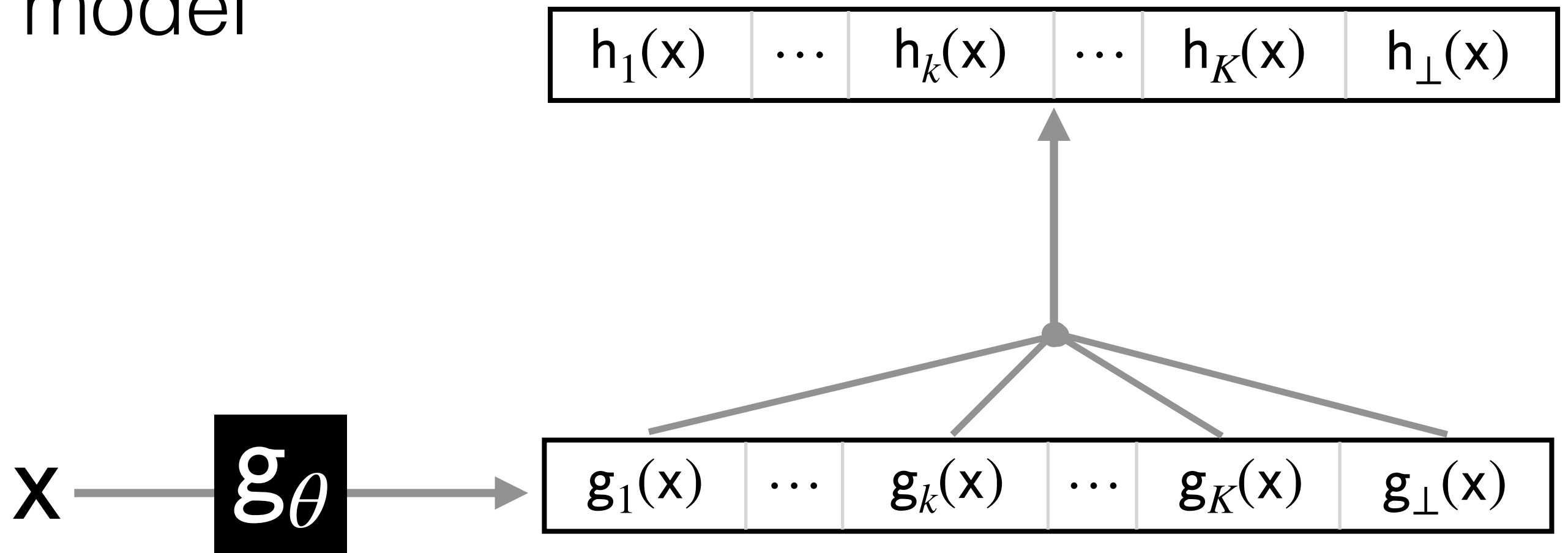
# softmax implementation

[Mozannar & Sontag, 2020]

training data

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

model



# softmax implementation

[Mozannar & Sontag, 2020]

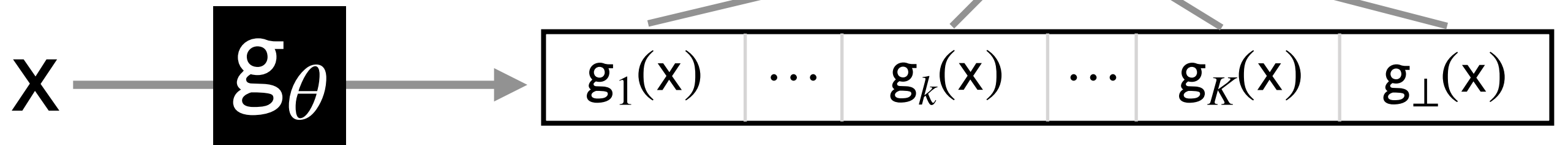
training data

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

model

$h_1(\mathbf{x})$	$\dots$	$h_k(\mathbf{x})$	$\dots$	$h_K(\mathbf{x})$	$h_{\perp}(\mathbf{x})$
-------------------	---------	-------------------	---------	-------------------	-------------------------

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



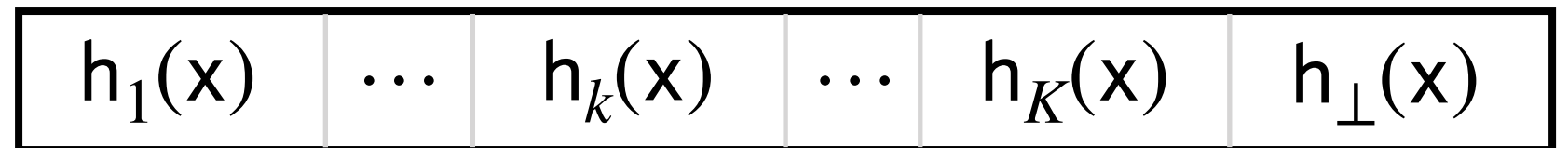
# softmax implementation

[Mozannar & Sontag, 2020]

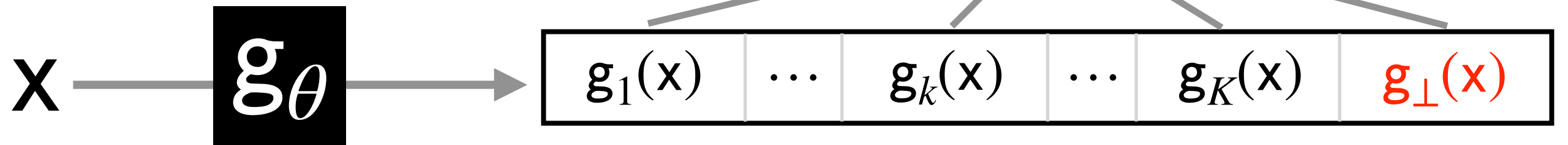
training data

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

model



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



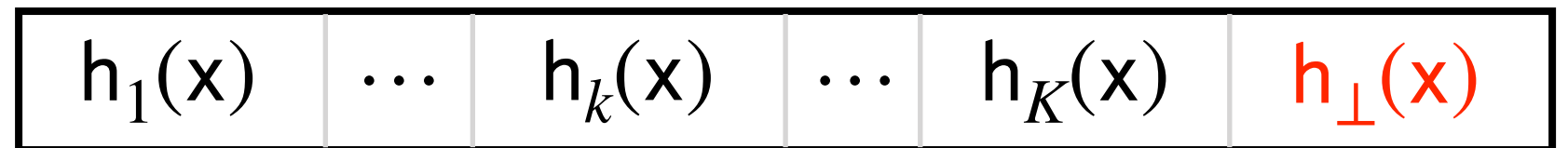
# softmax implementation

[Mozannar & Sontag, 2020]

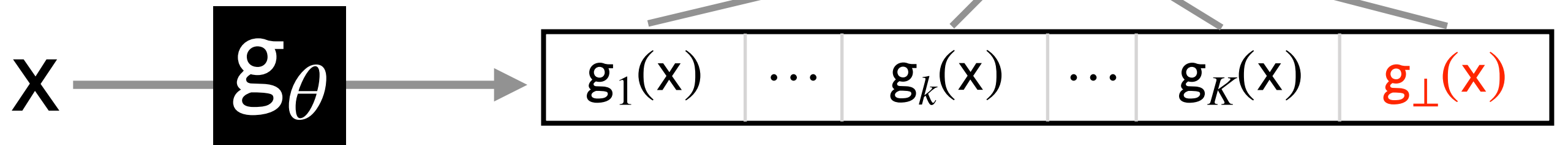
training data

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

model



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



# softmax implementation

[Mozannar & Sontag, 2020]

training data

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

model

$$h_i(\mathbf{x}) = \frac{\exp\{\mathbf{g}_i(\mathbf{x})\}}{\sum_{k=1}^{K+1} \exp\{\mathbf{g}_k(\mathbf{x})\}}$$

# softmax implementation

[Mozannar & Sontag, 2020]

training data

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

model

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

loss function

$$\ell(\theta; \mathbf{x}, y, \mathbf{m}) = -\log h_y(\mathbf{x}) - \mathbb{I}[y = \mathbf{m}] \cdot \log h_{\perp}(\mathbf{x})$$

# softmax implementation

[Mozannar & Sontag, 2020]

training data

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

model

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

loss function

$$\ell(\theta; \mathbf{x}, y, \mathbf{m}) = -\log h_y(\mathbf{x}) - \mathbb{I}[y = \mathbf{m}] \cdot \log h_{\perp}(\mathbf{x})$$



# softmax implementation

[Mozannar & Sontag, 2020]

training data

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

model

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

loss function

$$\ell(\theta; \mathbf{x}, y, \mathbf{m}) = -\log h_y(\mathbf{x}) - \mathbb{I}[y = \mathbf{m}] \cdot \log h_{\perp}(\mathbf{x})$$

# softmax implementation

[Mozannar & Sontag, 2020]

training data

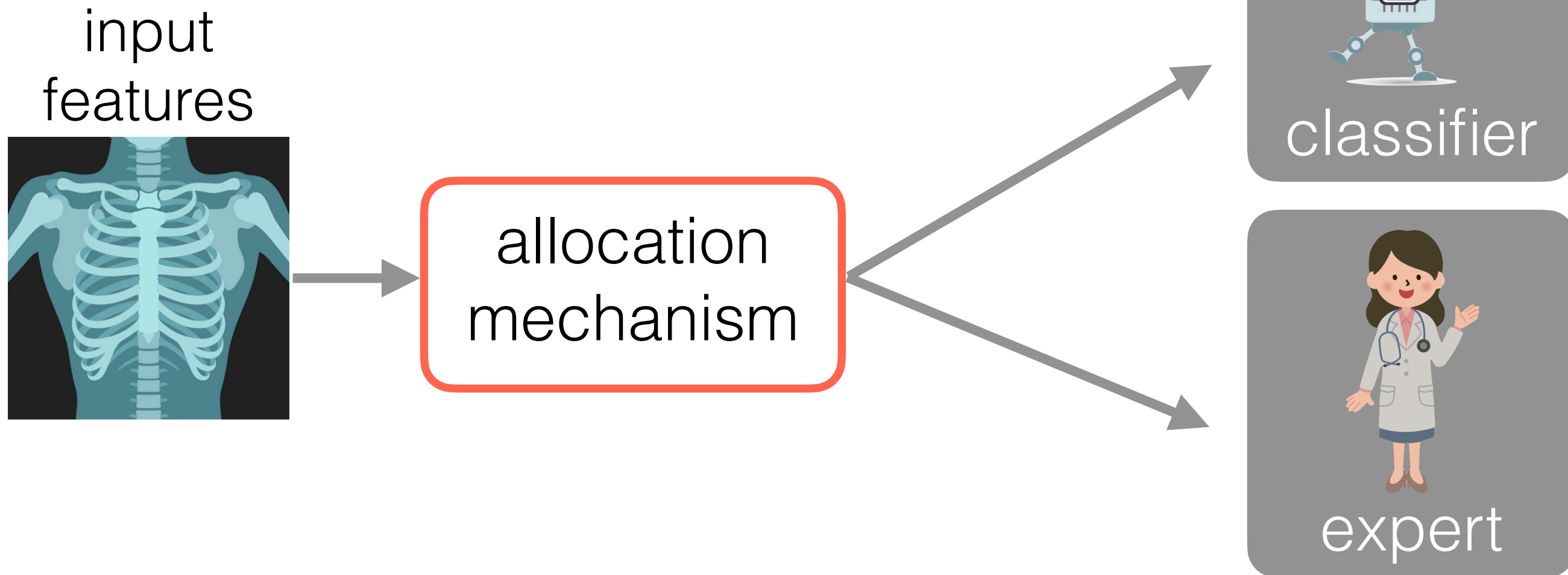
$$\mathfrak{D} = \left\{ \mathbf{x}_n, y_n, \mathbf{m}_n \right\}_{n=1}^N$$

model

$$h_i(\mathbf{x}) = \frac{\exp\{\mathbf{g}_i(\mathbf{x})\}}{\sum_{k=1}^{K+1} \exp\{\mathbf{g}_k(\mathbf{x})\}}$$

loss function

$$\ell(\theta; \mathbf{x}, y, \mathbf{m}) = -\log h_y(\mathbf{x}) - \mathbb{I}[y = \mathbf{m}] \cdot \log h_{\perp}(\mathbf{x})$$



defer to expert if...

$$\max_{y \in [1, K]} h_y(\mathbf{x}) \leq h_{\perp}(\mathbf{x})$$

## ⊗ single expert

- ⊗ softmax surrogate loss

- ⊗ improving calibration via one-vs-all

## ⊗ multiple experts

- ⊗ surrogate losses

- ⊗ conformal sets of experts

## ⊗ population of experts

- ⊗ surrogate losses

- ⊗ meta-learning a rejector

## ⊗ single expert

- ⊗ softmax surrogate loss

- ⊗ improving calibration via one-vs-all

## ⊗ multiple experts

- ⊗ surrogate losses

- ⊗ conformal sets of experts

## ⊗ population of experts

- ⊗ surrogate losses

- ⊗ meta-learning a rejector

How well can the softmax-based system estimate expert correctness?

$$\hat{p}(m = y | x) \stackrel{?}{\approx} \mathbb{P}(m = y | x)$$

How well can the softmax-based system estimate expert correctness?

$$\hat{p}(m = y | x) \underset{?}{\approx} \mathbb{P}(m = y | x)$$

- ⊗ optimal allocation
- ⊗ transparency
- ⊗ detecting distribution shift  
(in the expert)

How well can the softmax-based system estimate expert correctness?

$$\hat{p}(m = y | x) \not\approx \mathbb{P}(m = y | x)$$



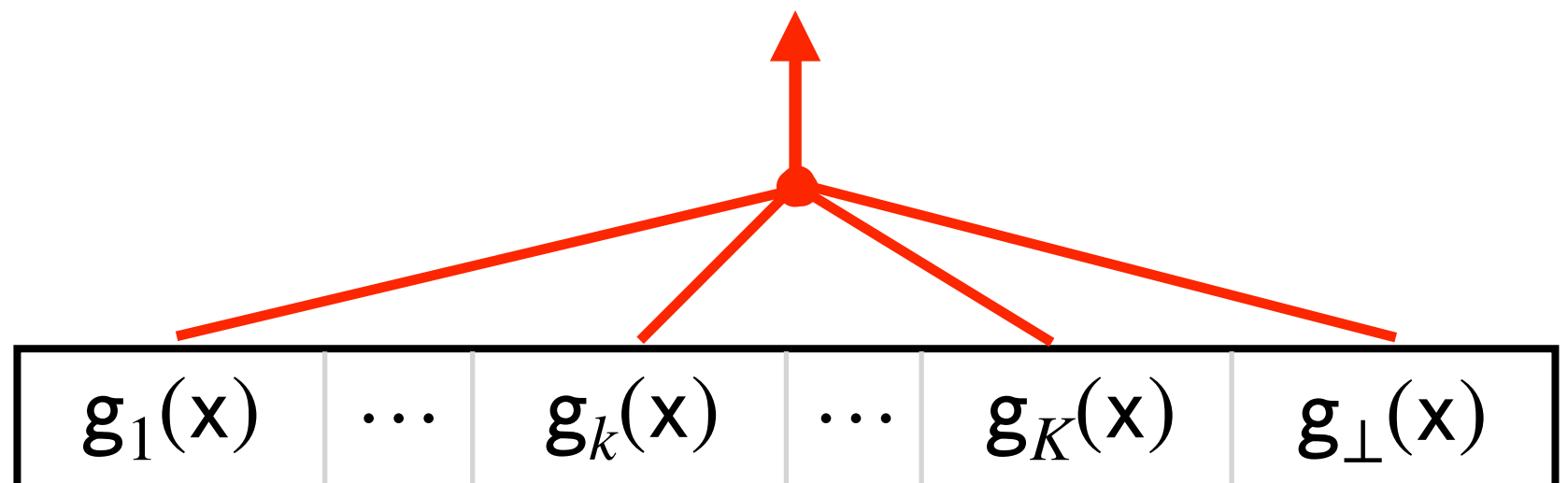
How well can the softmax-based system estimate expert correctness?

$$\hat{p}(m = y | x) \not\approx \mathbb{P}(m = y | x)$$

degenerate  
parameterization

[Proposition 3.1]

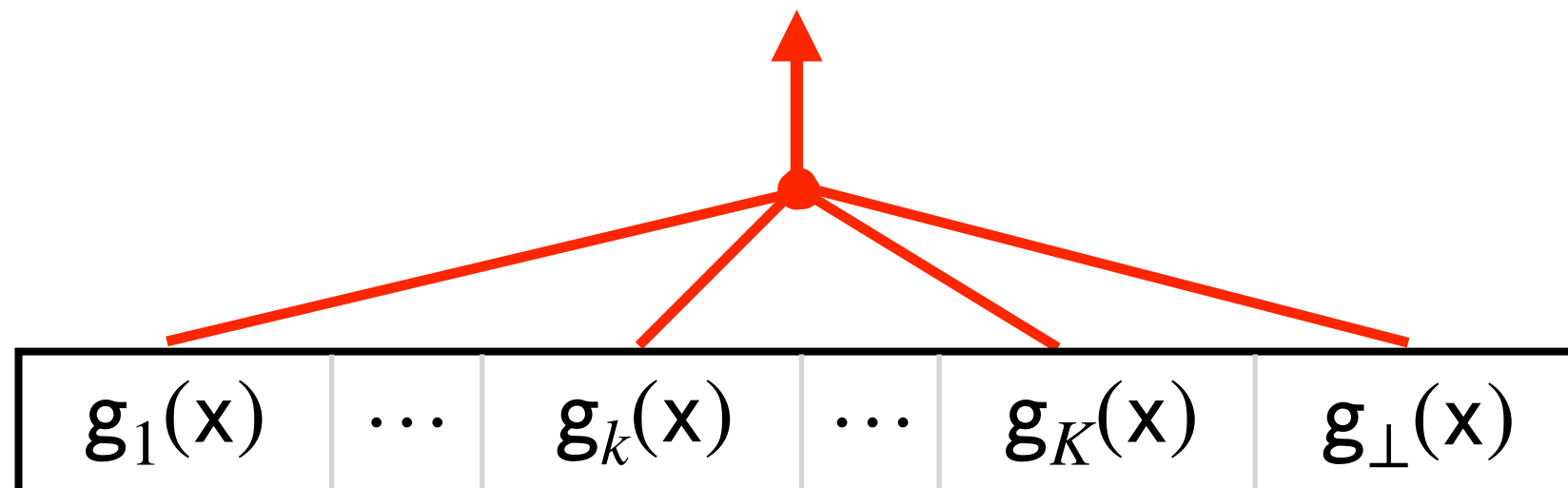
$$h_{\perp}(x) = \frac{\exp\{g_{\perp}(x)\}}{\sum_{k=1}^{K+1} \exp\{g_k(x)\}}$$



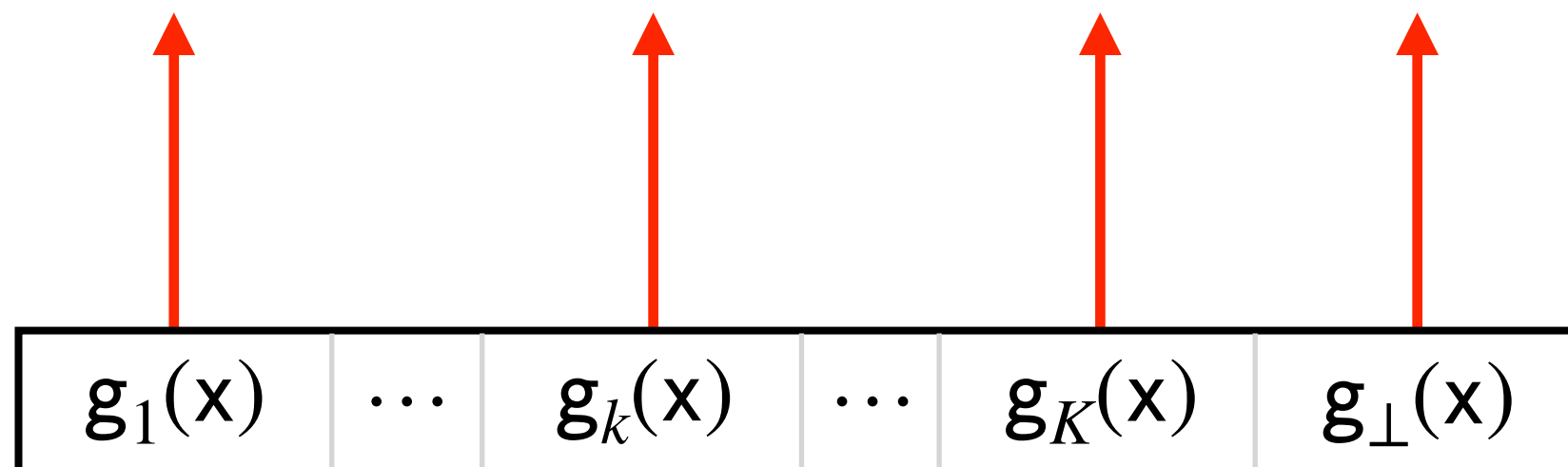
our solution: switch to a  
one-vs-all parameterization

**our solution:** switch to a  
one-vs-all parameterization

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

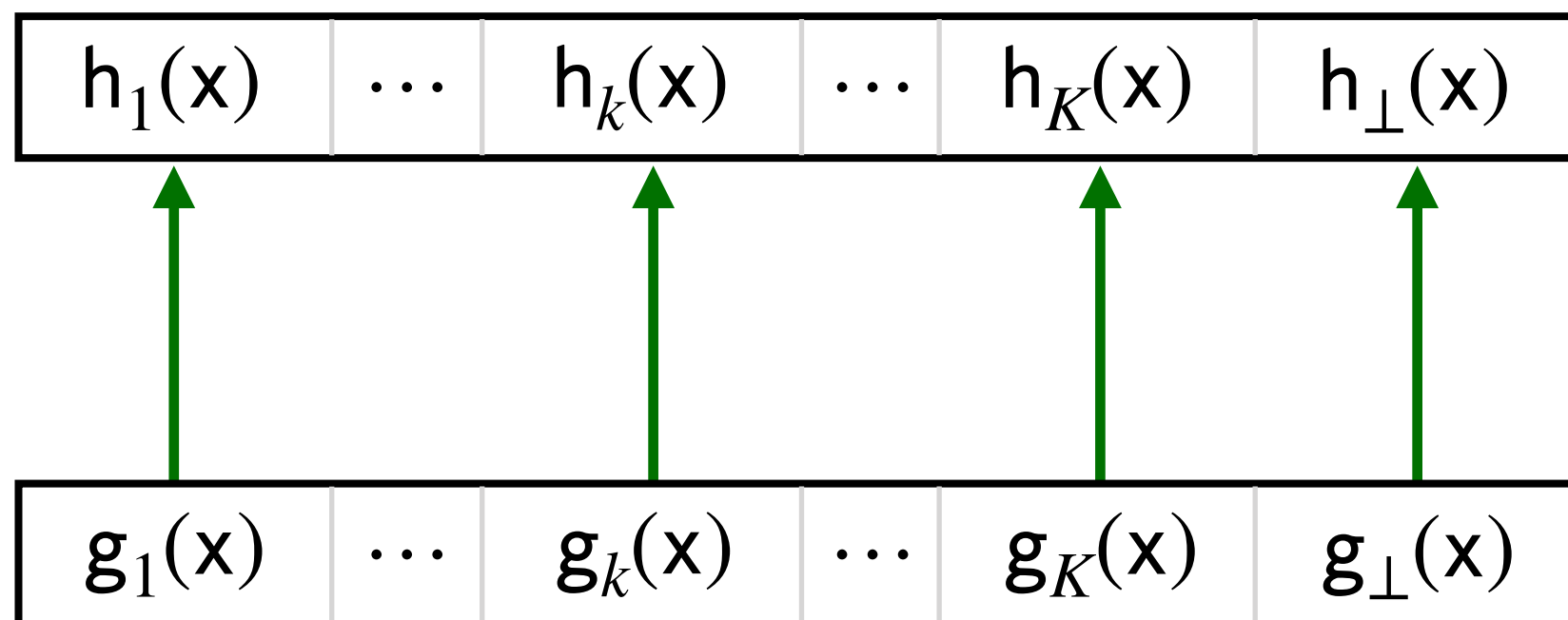


**our solution:** switch to a  
one-vs-all parameterization



**our solution:** switch to a  
one-vs-all parameterization

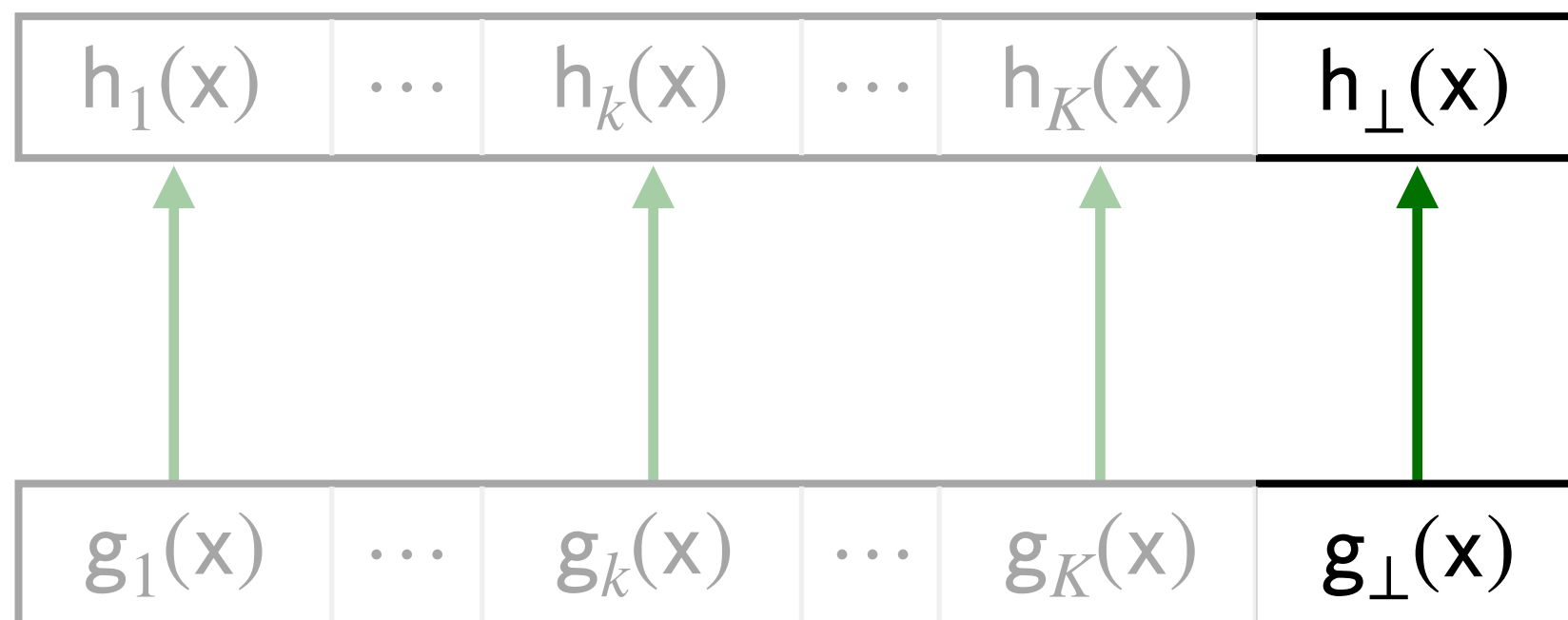
$$h_i(\mathbf{x}) = \frac{1}{1 + \exp \{ -g_i(\mathbf{x}) \}}$$



**our solution:** switch to a  
one-vs-all parameterization

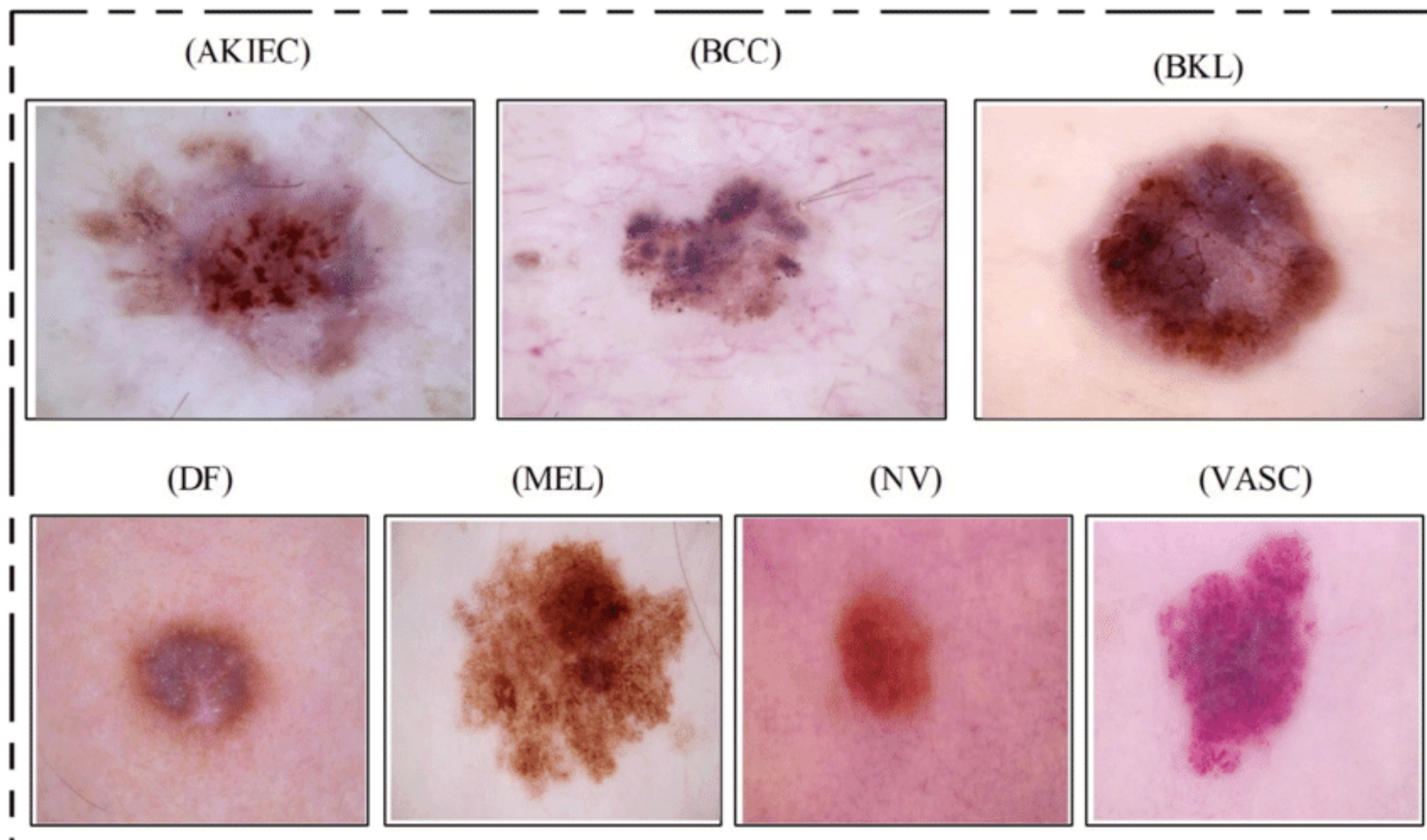
[Theorem 4.1]

$$h_{\perp}^*(\mathbf{x}) = \mathbb{P}(m = y \mid \mathbf{x})$$



# estimating expert correctness

## skin lesion diagnosis



# estimating expert correctness

$\hat{p}$

distance:  $\hat{p}$  vs  $\mathbb{P}$

---

softmax

one-vs-all  
(ours)

---



# estimating expert correctness

$\hat{p}$

distance:  $\hat{p}$  vs  $\mathbb{P}$

---

softmax

$26.7 \pm 1.8$

---

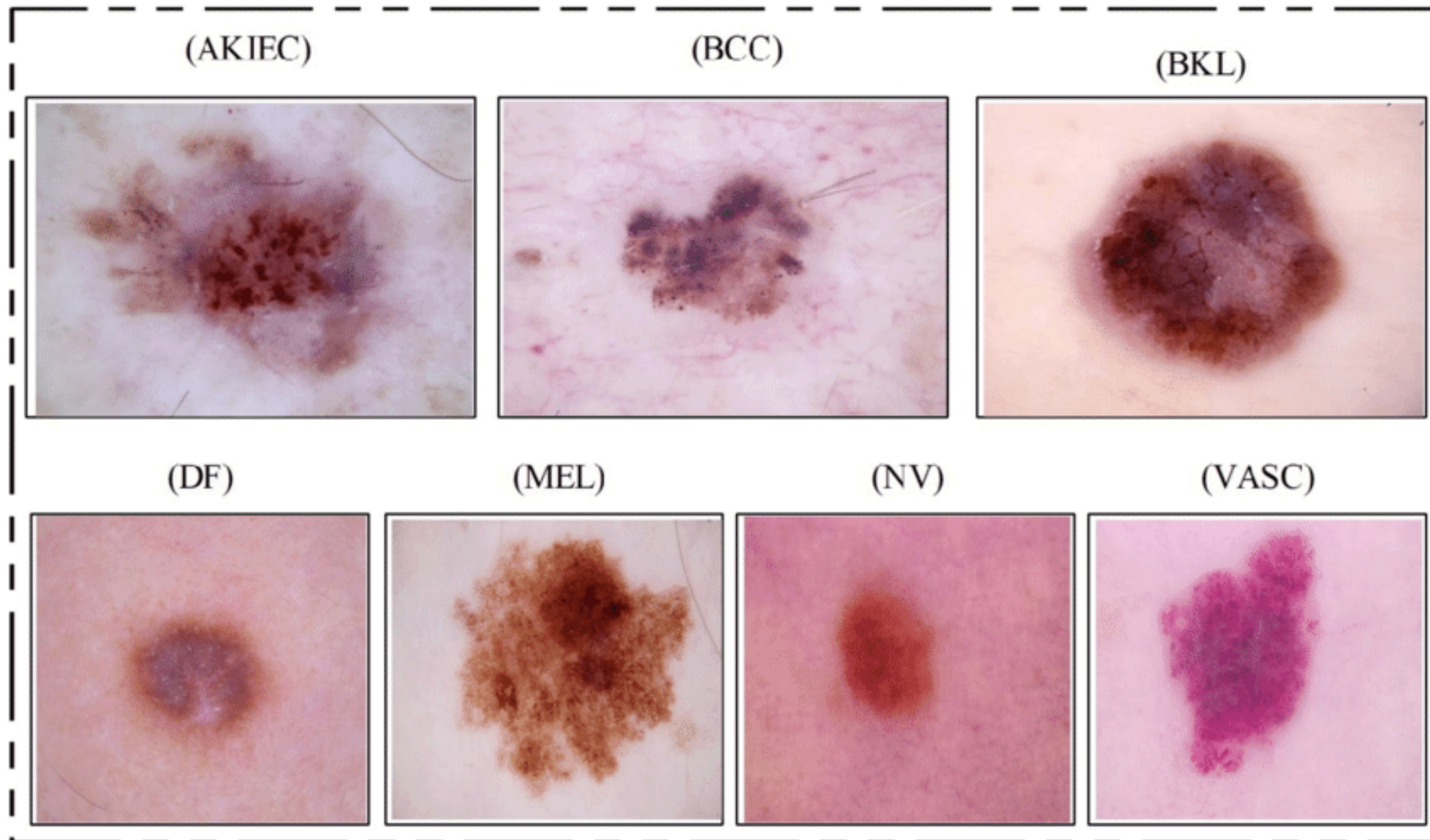
one-vs-all  
(ours)

$8.0 \pm 1.0$



But does one-vs-all  
result in more accurate models?

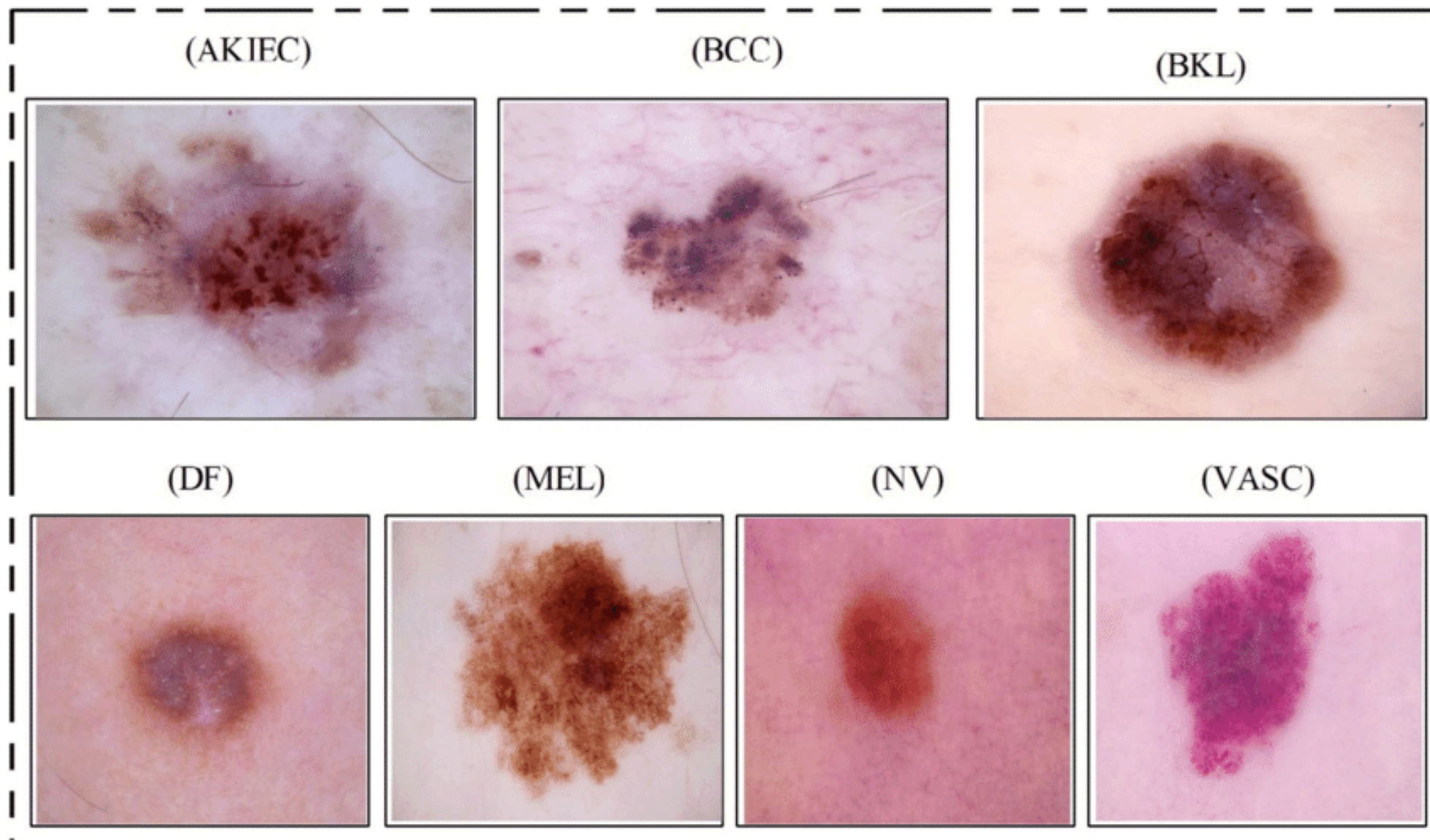
# skin lesion diagnosis




—●— one-vs-all (ours)

—●— softmax

# skin lesion diagnosis



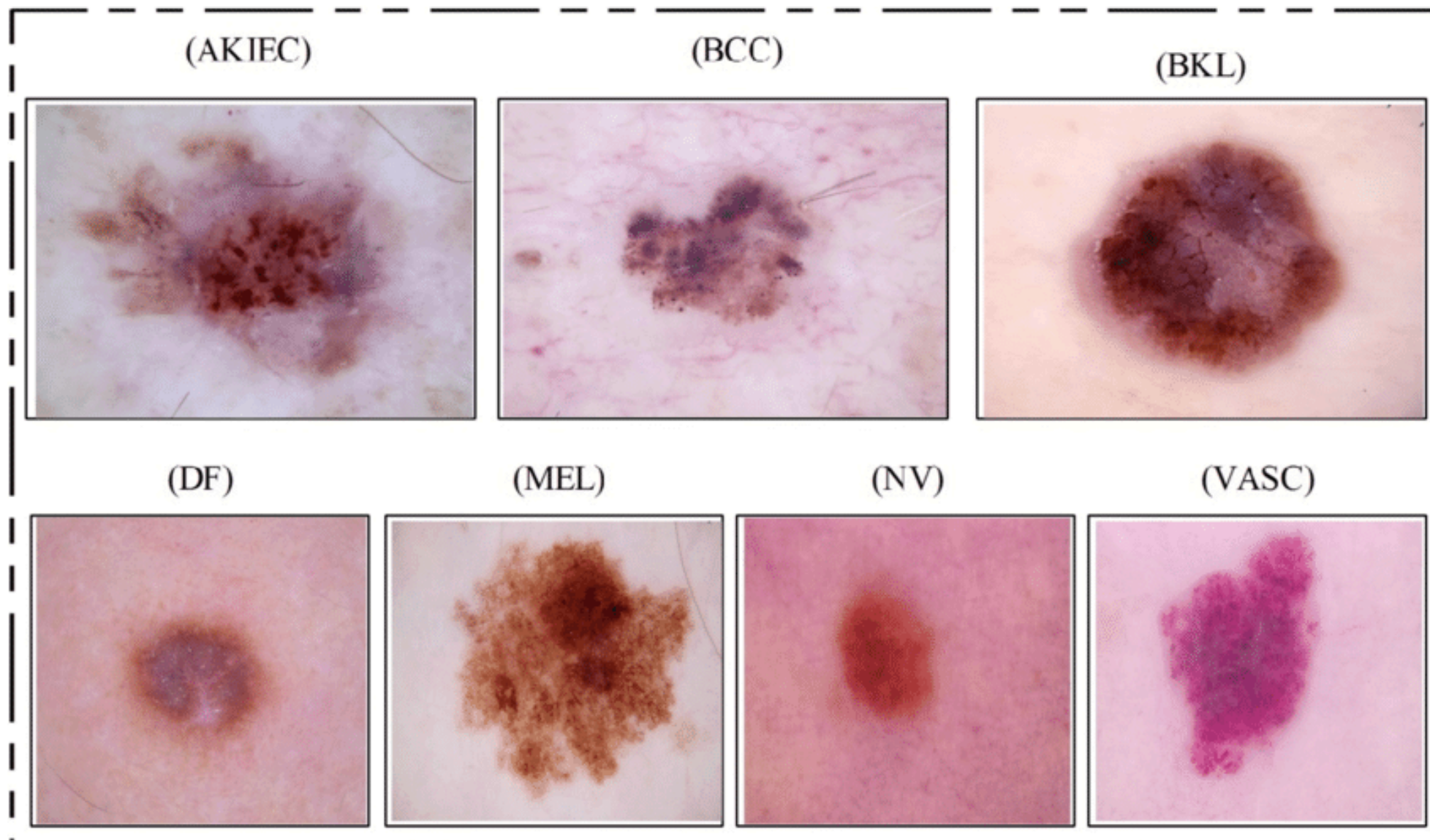


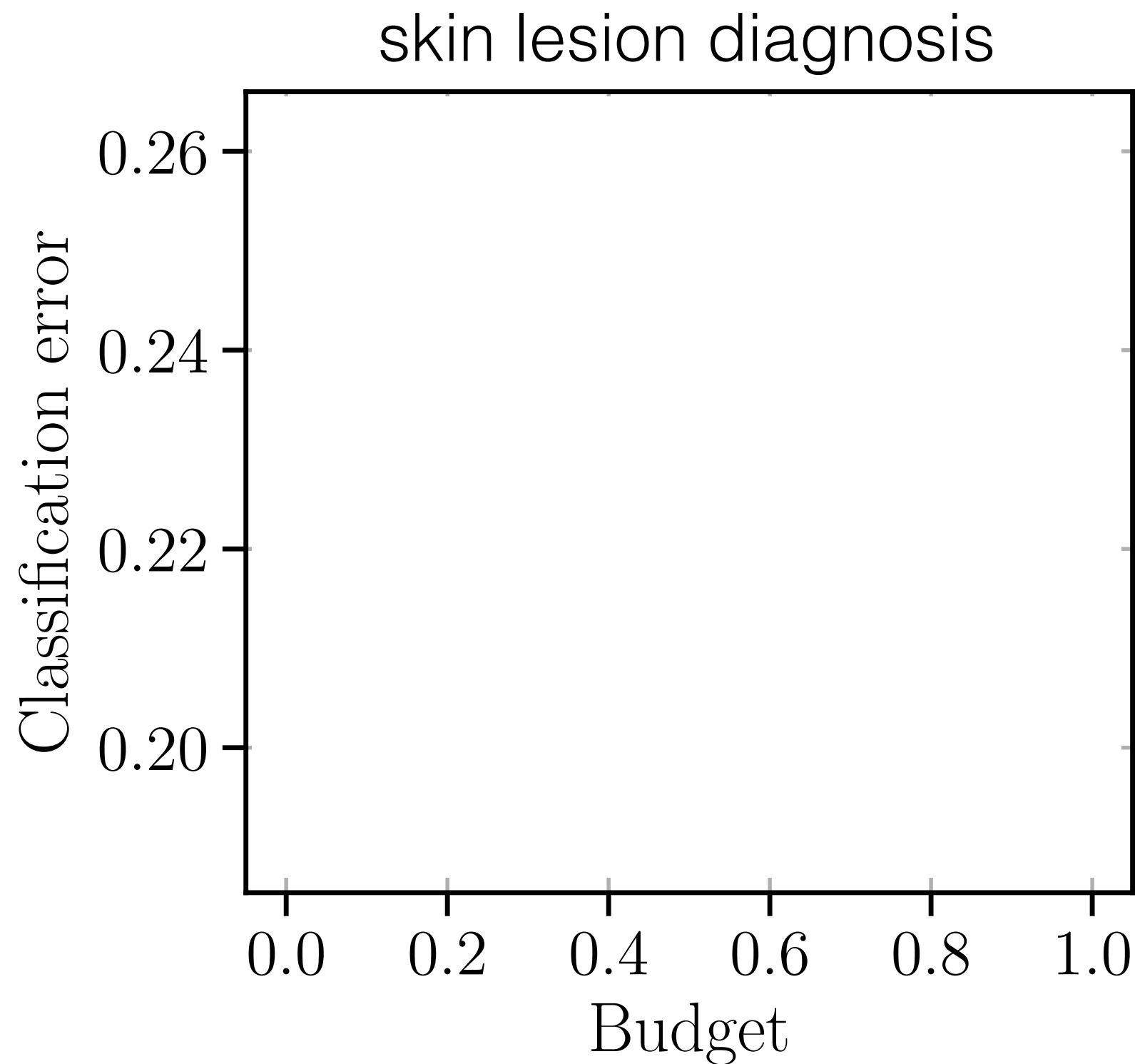
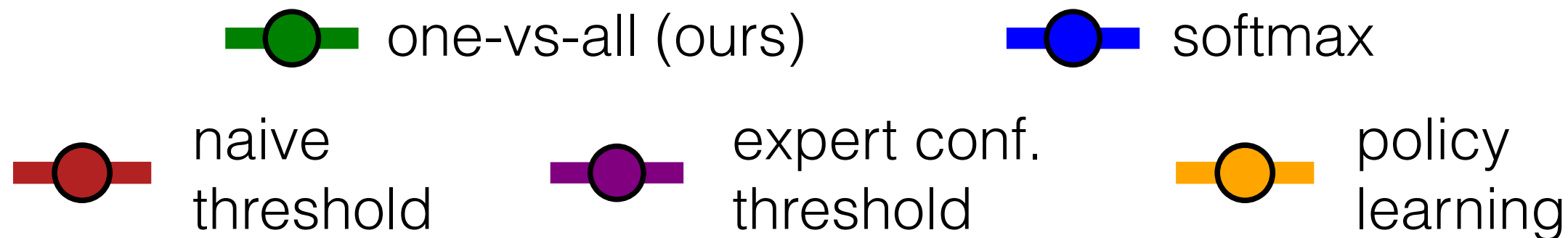


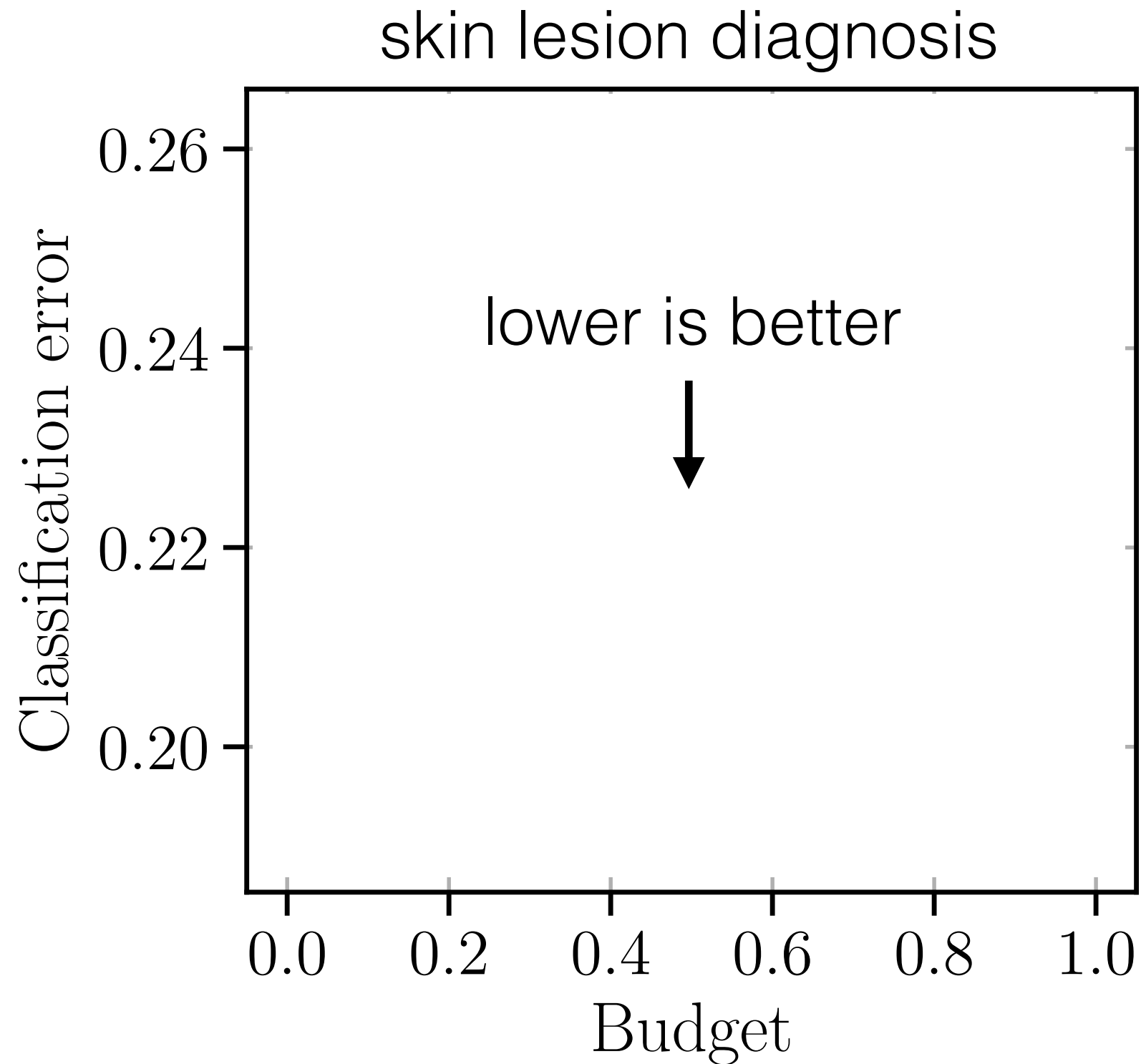
Legend for the methods compared in the paper:

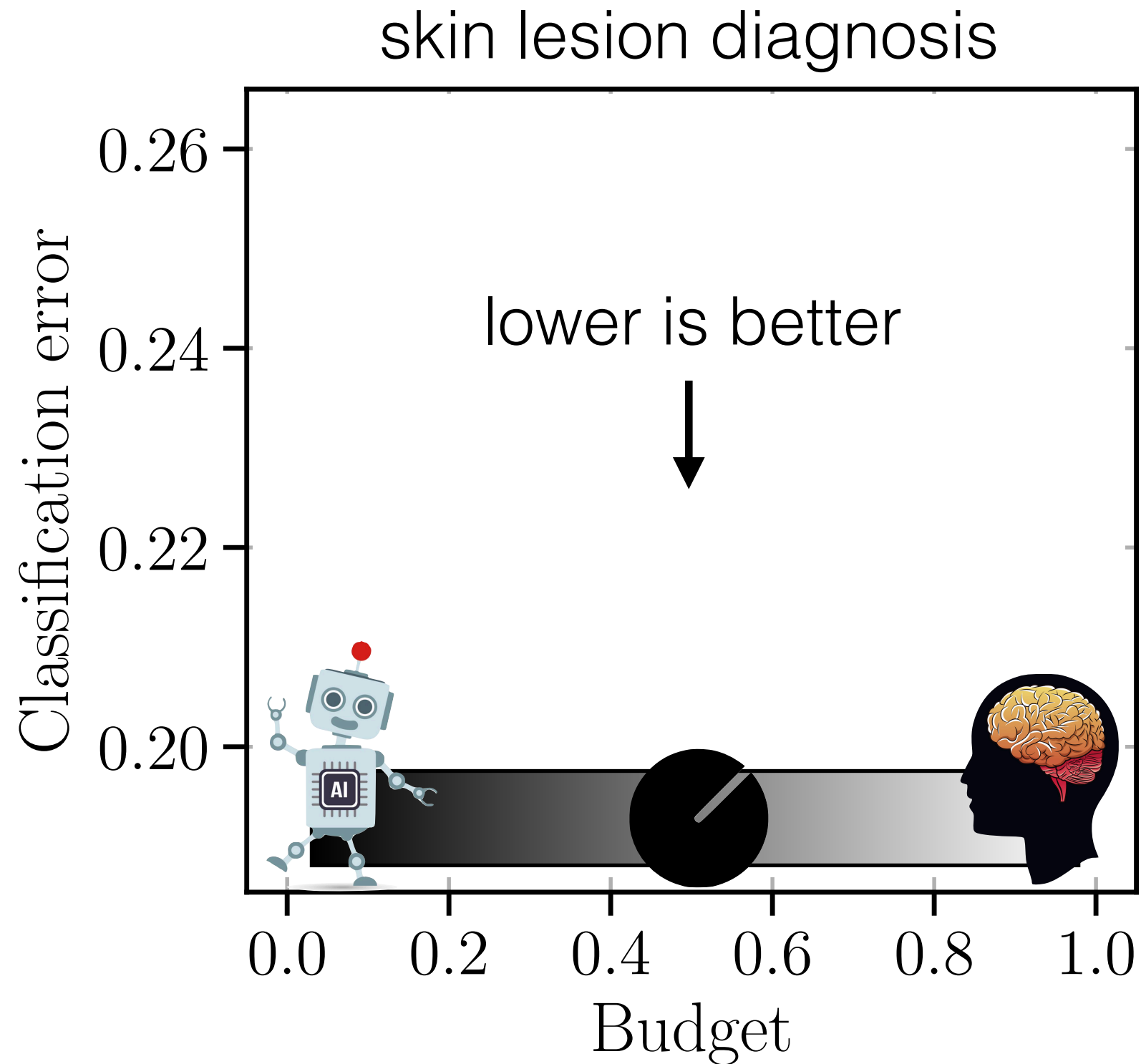
- one-vs-all (ours)
- softmax
- naive threshold
- expert conf. threshold
- policy learning

# skin lesion diagnosis

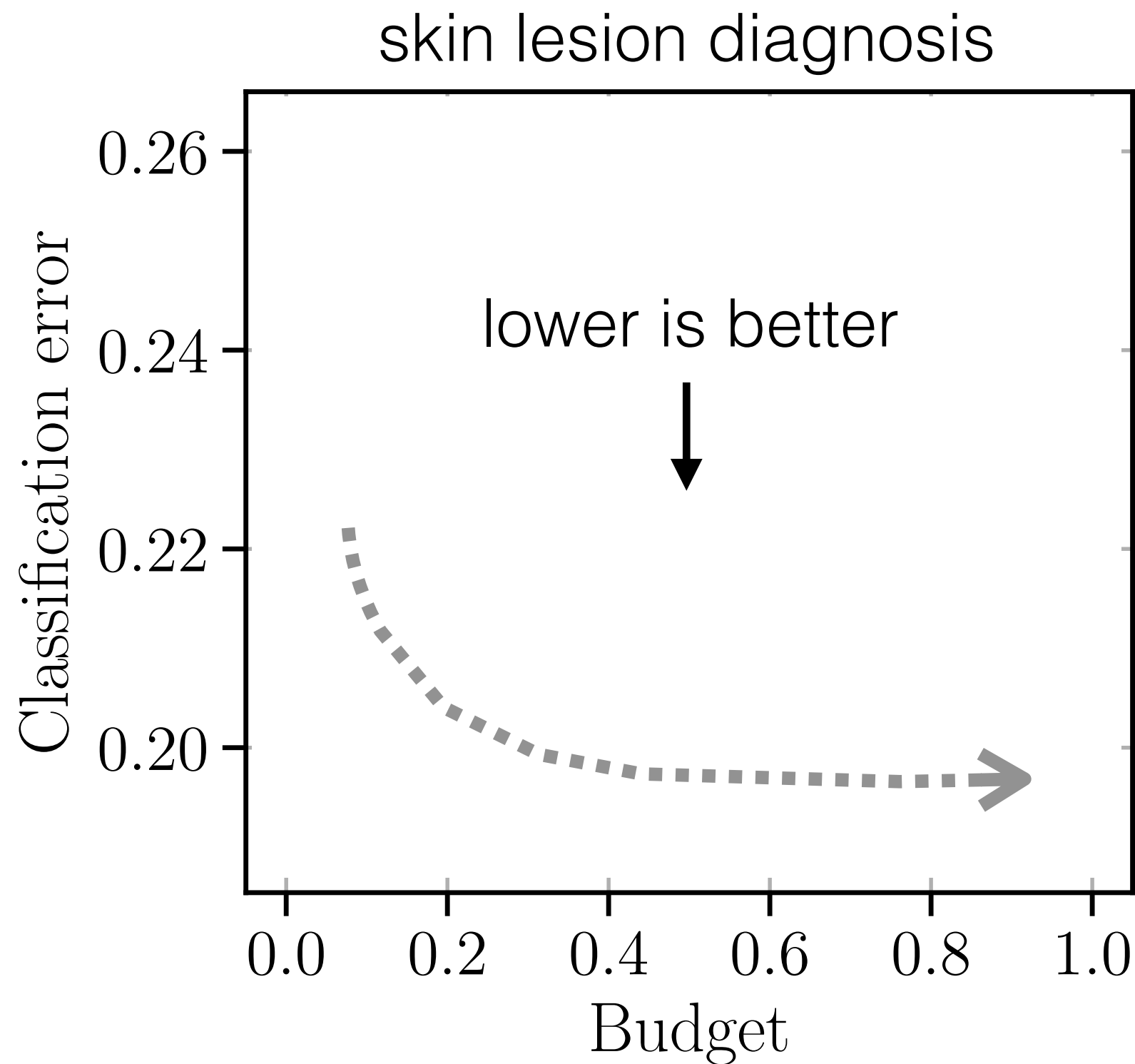
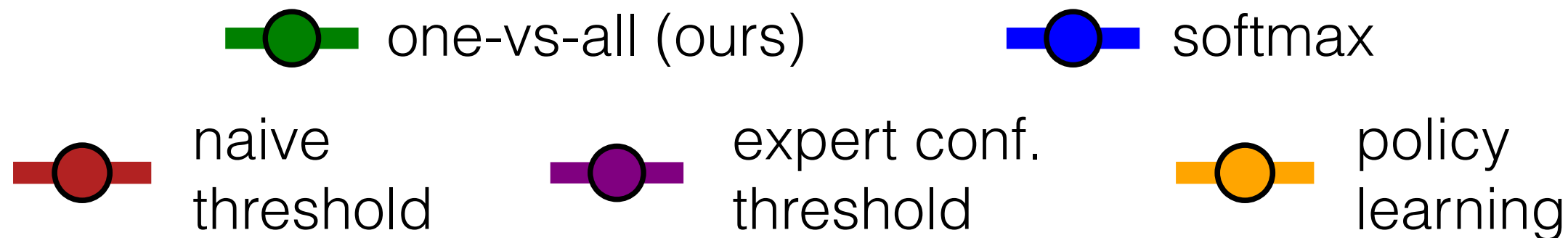


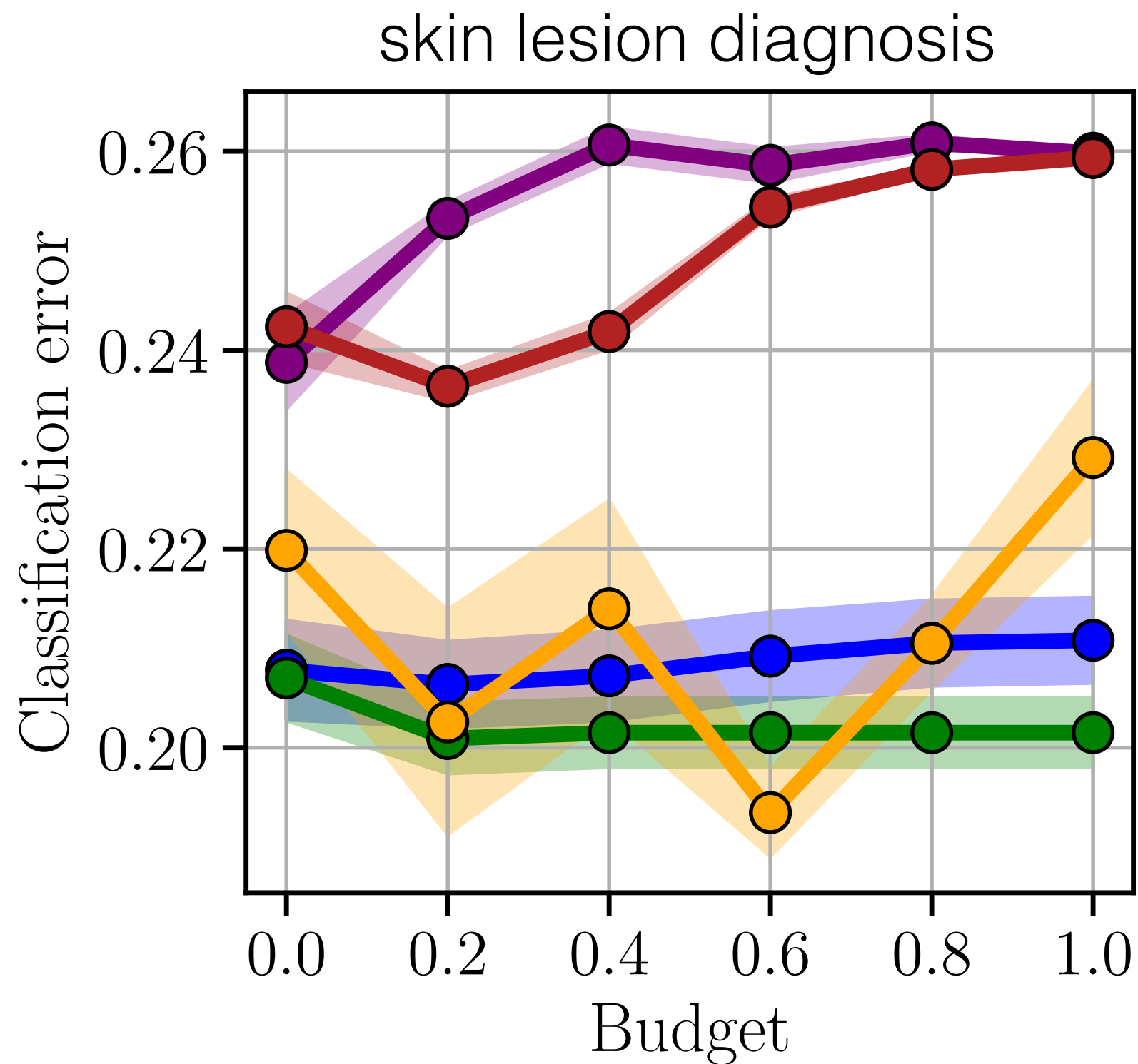
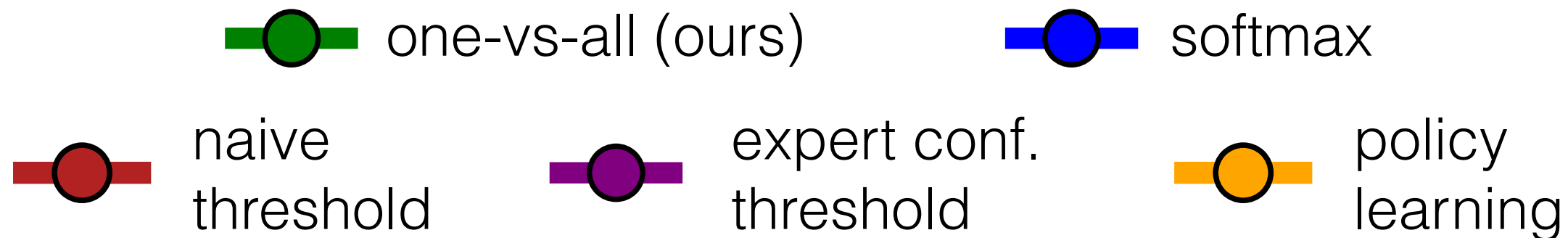


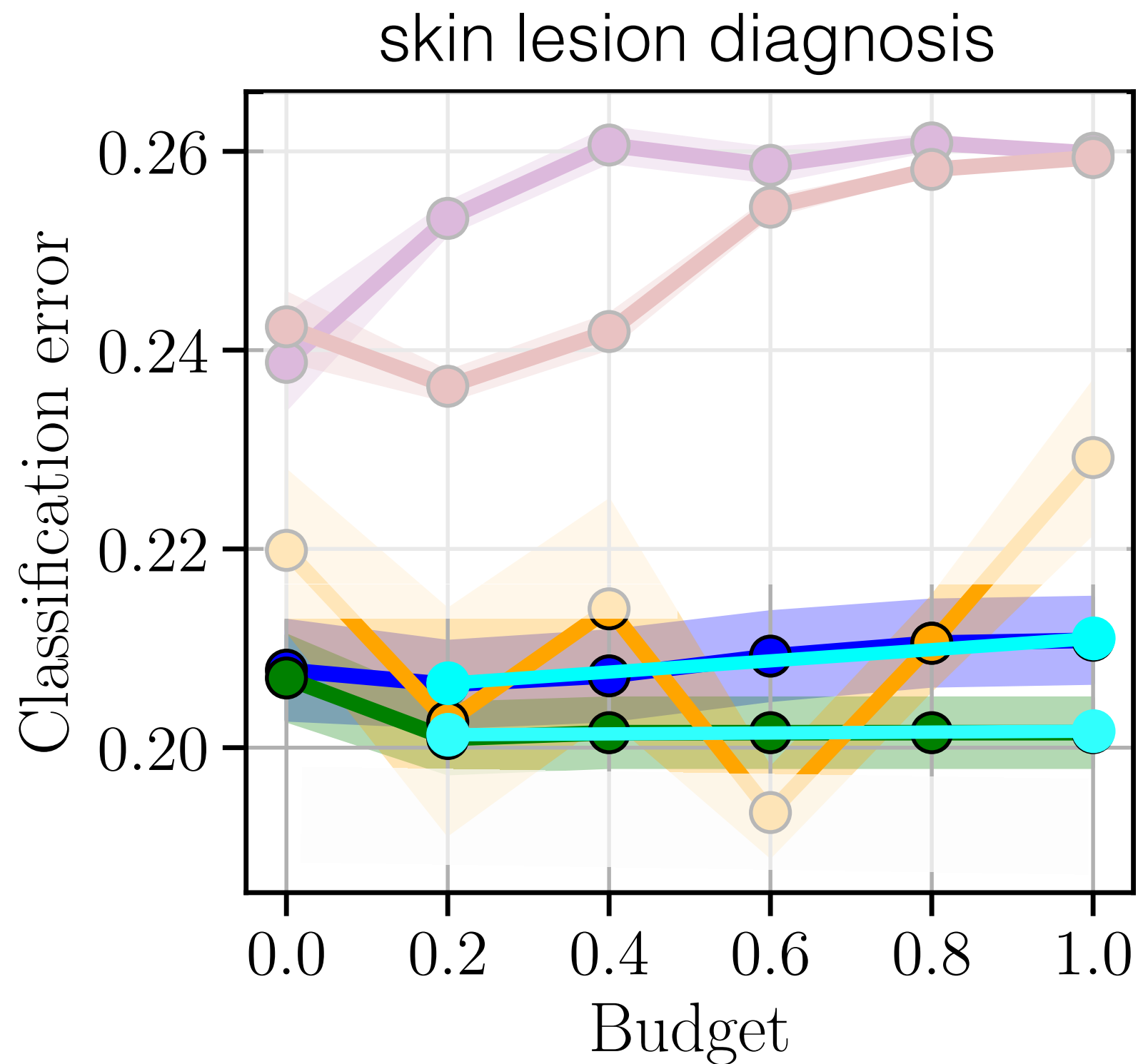
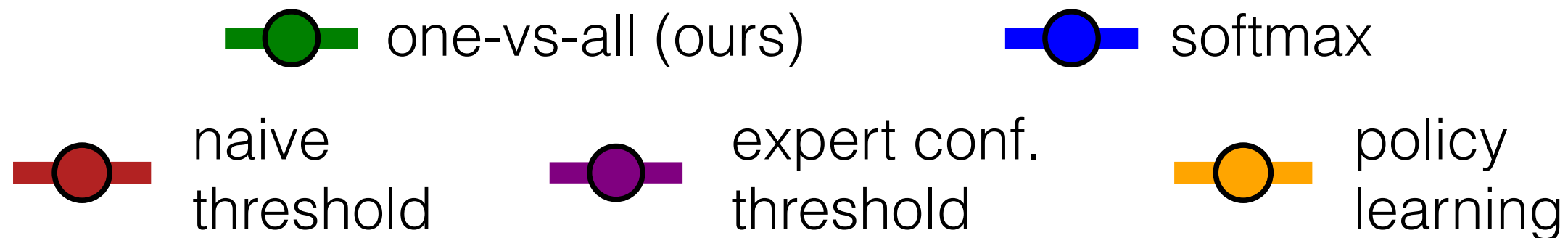












## ⊗ single expert

- ⊗ softmax surrogate loss

- ⊗ improving calibration via one-vs-all

## ⊗ multiple experts

- ⊗ surrogate losses

- ⊗ conformal sets of experts

## ⊗ population of experts

- ⊗ surrogate losses

- ⊗ meta-learning a rejector

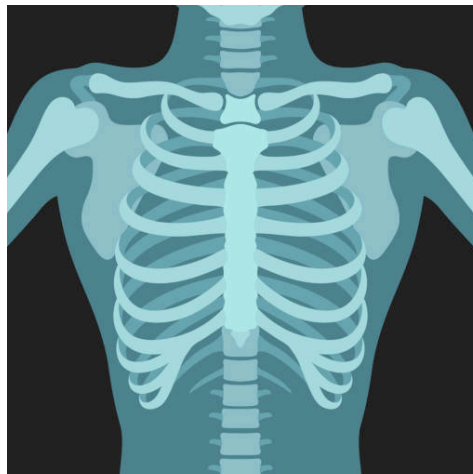
- ⊗ single expert
  - ⊗ softmax surrogate loss
  - ⊗ improving calibration via one-vs-all

- ⊗ **multiple experts**

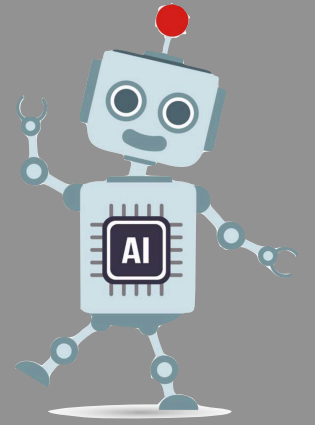
- ⊗ surrogate losses
  - ⊗ conformal sets of experts

- ⊗ population of experts
  - ⊗ surrogate losses
  - ⊗ meta-learning a rejector

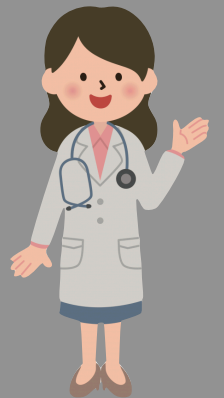
input  
features



allocation  
mechanism

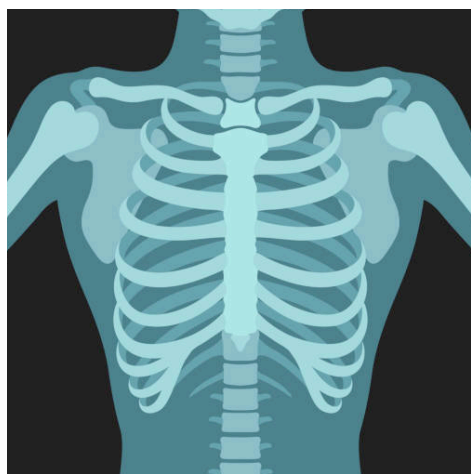


classifier

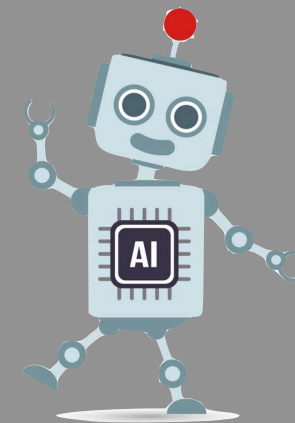


expert

input  
features



allocation  
mechanism



classifier



expert #1

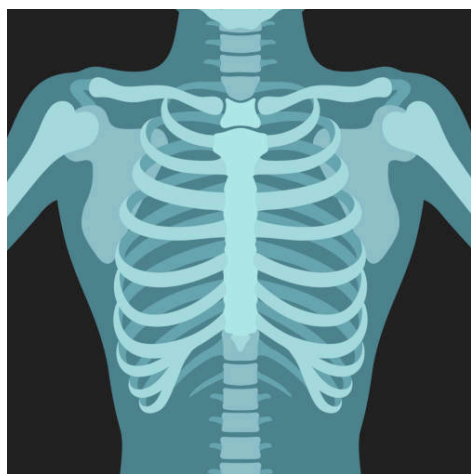


expert #3

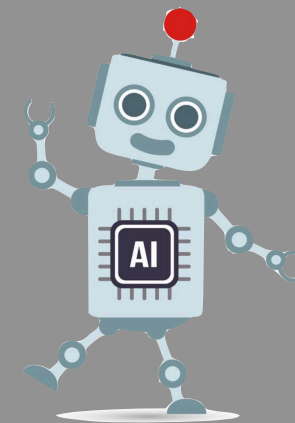


expert #2

input  
features



allocation  
mechanism



classifier



expert #1



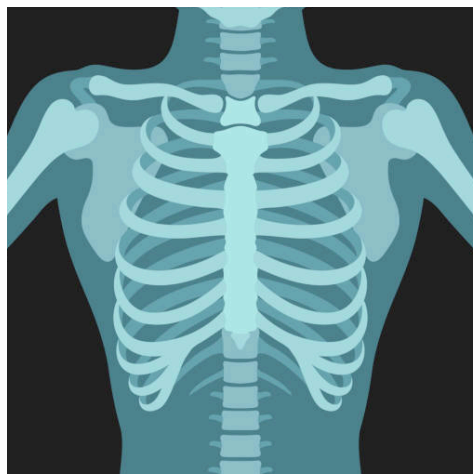
expert #3



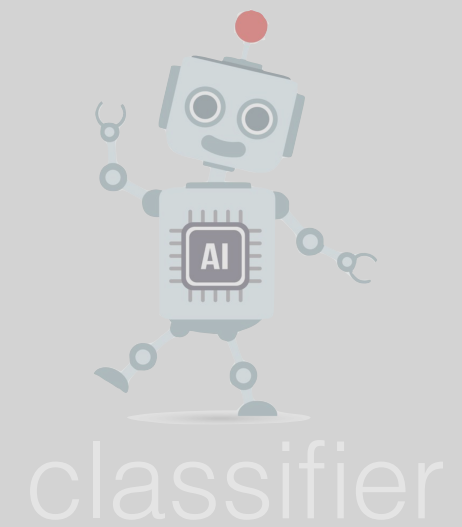
expert #2



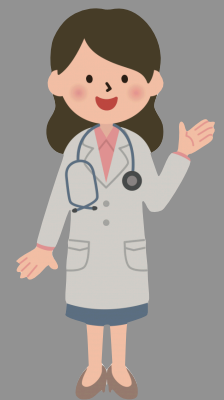
input  
features



allocation  
mechanism



classifier



expert #1

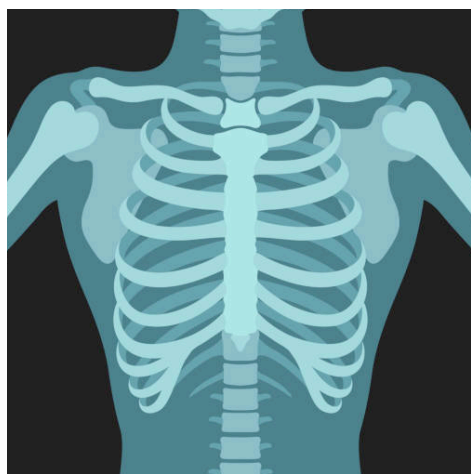


expert #3

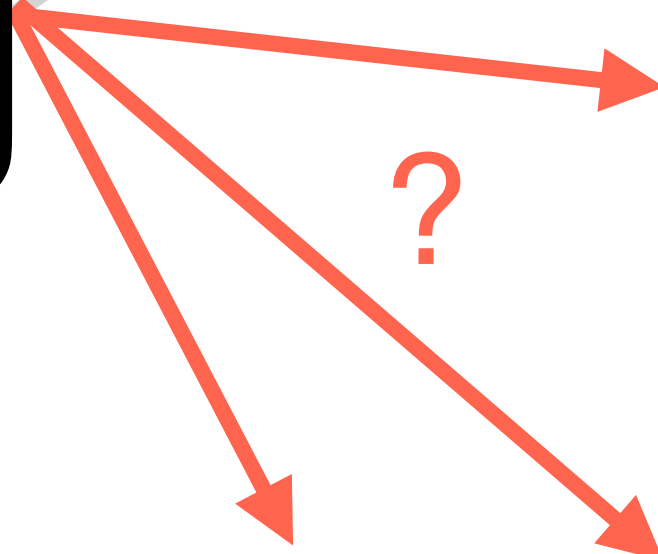
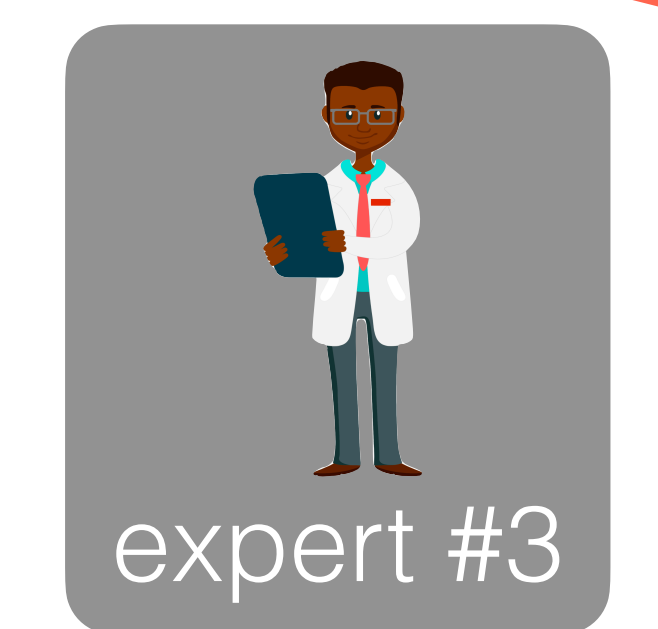
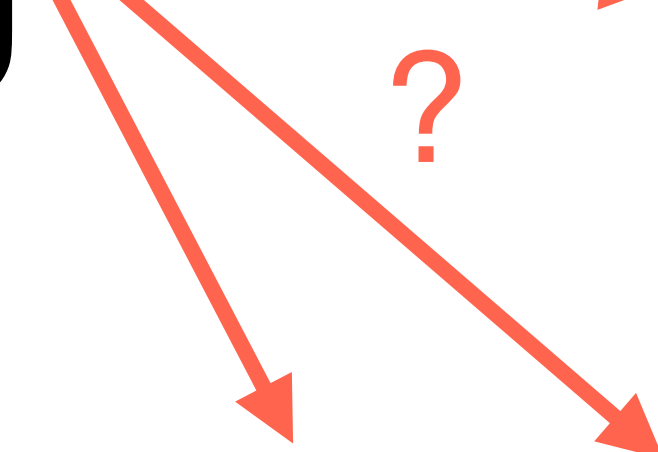


expert #2

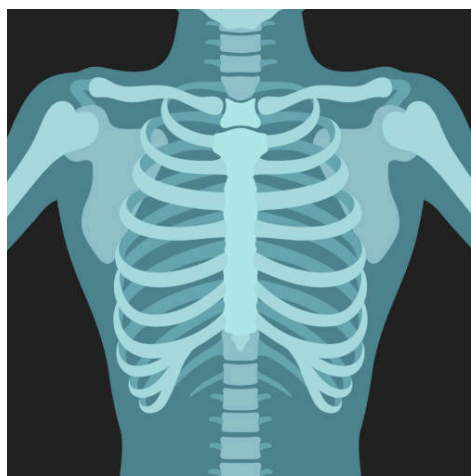
input  
features



allocation  
mechanism



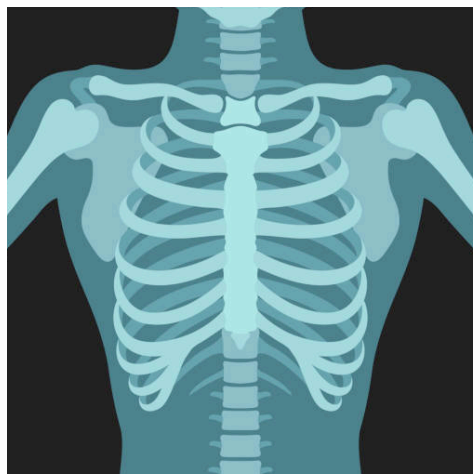
input  
features



allocation  
mechanism



input  
features



allocation  
mechanism



$$\mathbb{P}(m_1 = y | x)$$

expert #1

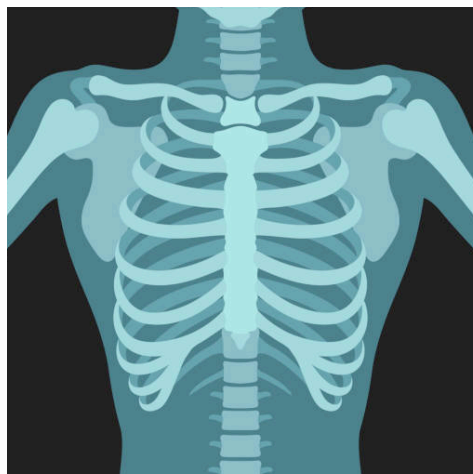
$$\mathbb{P}(m_3 = y | x)$$

expert #3

$$\mathbb{P}(m_2 = y | x)$$

expert #2

input  
features



allocation  
mechanism

$\mathbb{P}(y | \mathbf{x})$

classifier

$\mathbb{P}(m_1 = y | \mathbf{x})$

expert #1

$\mathbb{P}(m_3 = y | \mathbf{x})$

expert #3

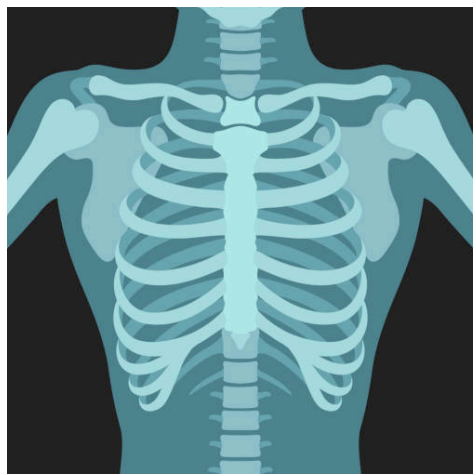
$\mathbb{P}(m_2 = y | \mathbf{x})$

expert #2

use classifier if...

$$\max_y \mathbb{P}(y | \mathbf{x}) > \mathbb{P}(m_j = y | \mathbf{x}), \forall j$$

input  
features



allocation  
mechanism

$$\mathbb{P}(y | \mathbf{x})$$

classifier

$$\mathbb{P}(m_1 = y | \mathbf{x})$$

expert #1

$$\mathbb{P}(m_3 = y | \mathbf{x})$$

expert #3

$$\mathbb{P}(m_2 = y | \mathbf{x})$$

expert #2

else, pick best expert:

$$\arg \max_j \mathbb{P}(m_j = y | \mathbf{x})$$

# multi-expert implementation

training data

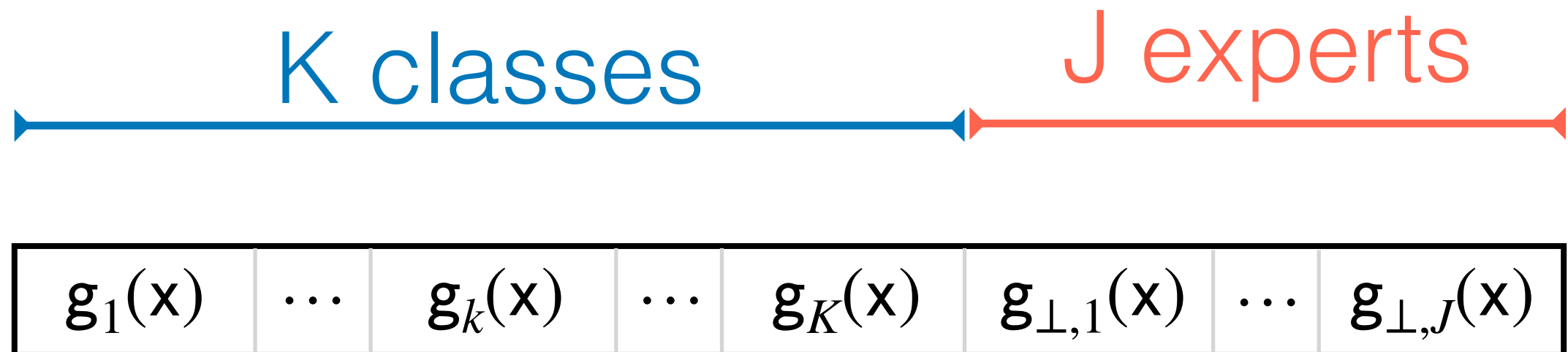
$$\mathfrak{D} = \left\{ \mathbf{x}_n, y_n, \mathbf{m}_{n,1}, \dots, \mathbf{m}_{n,J} \right\}_{n=1}^N$$

# multi-expert implementation

training data

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

model



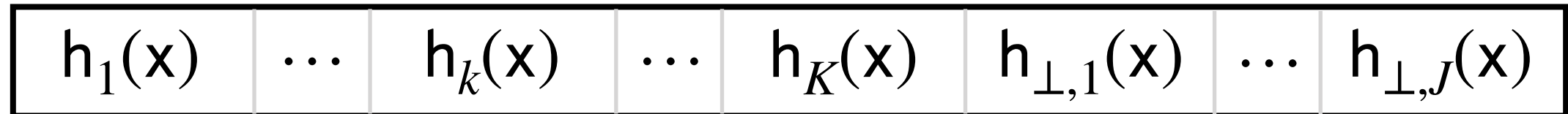


# multi-expert implementation

training data

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

model



K classes

J experts

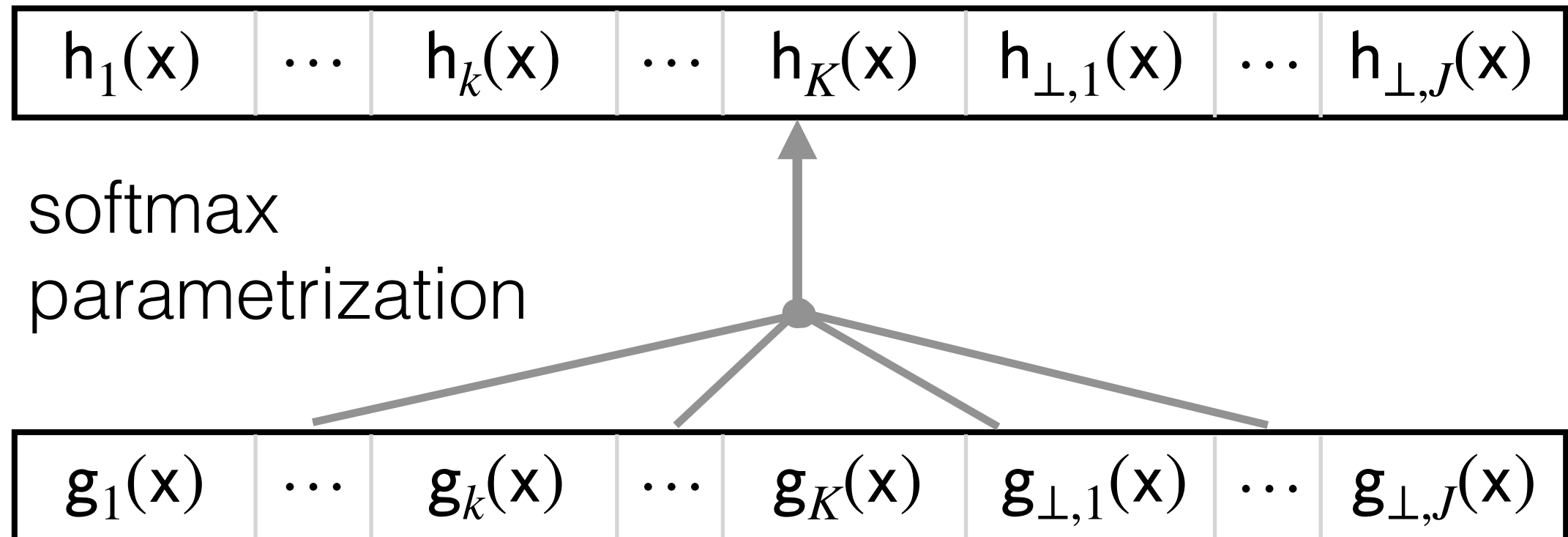


# multi-expert implementation

training data

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

model

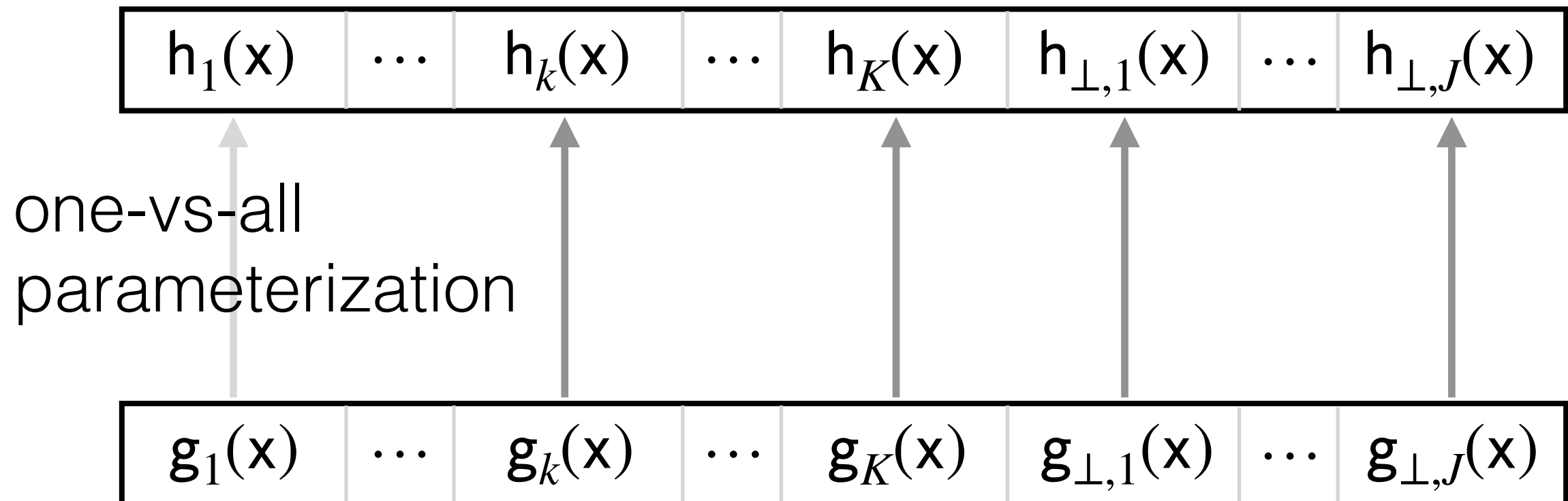


# multi-expert implementation

training data

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

model



# multi-expert implementation

training data

$$\mathfrak{D} = \left\{ \mathbf{x}_n, y_n, \mathbf{m}_{n,1}, \dots, \mathbf{m}_{n,J} \right\}_{n=1}^N$$

model

- ⊗ softmax and one-vs-all variants
- ⊗ both consistent w.r.t. 0-1 loss

# multi-expert implementation

training data

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

model

- ⊗ softmax and one-vs-all variants
- ⊗ both consistent w.r.t. 0-1 loss

softmax loss function

$$\ell(\theta; \mathbf{x}, y, \mathbf{m}) = -\log h_y(\mathbf{x}) - \sum_j \mathbb{I}[y = m_j] \cdot \log h_{\perp,j}(\mathbf{x})$$

# multi-expert implementation

training data

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

model

- ⊗ softmax and one-vs-all variants
- ⊗ both consistent w.r.t. 0-1 loss

softmax loss function

$$\ell(\theta; \mathbf{x}, y, \mathbf{m}) = -\log h_y(\mathbf{x}) - \sum_j \mathbb{I}[y = \mathbf{m}_j] \cdot \log h_{\perp,j}(\mathbf{x})$$

# multi-expert implementation

training data

$$\mathfrak{D} = \left\{ \mathbf{x}_n, y_n, \mathbf{m}_{n,1}, \dots, \mathbf{m}_{n,J} \right\}_{n=1}^N$$

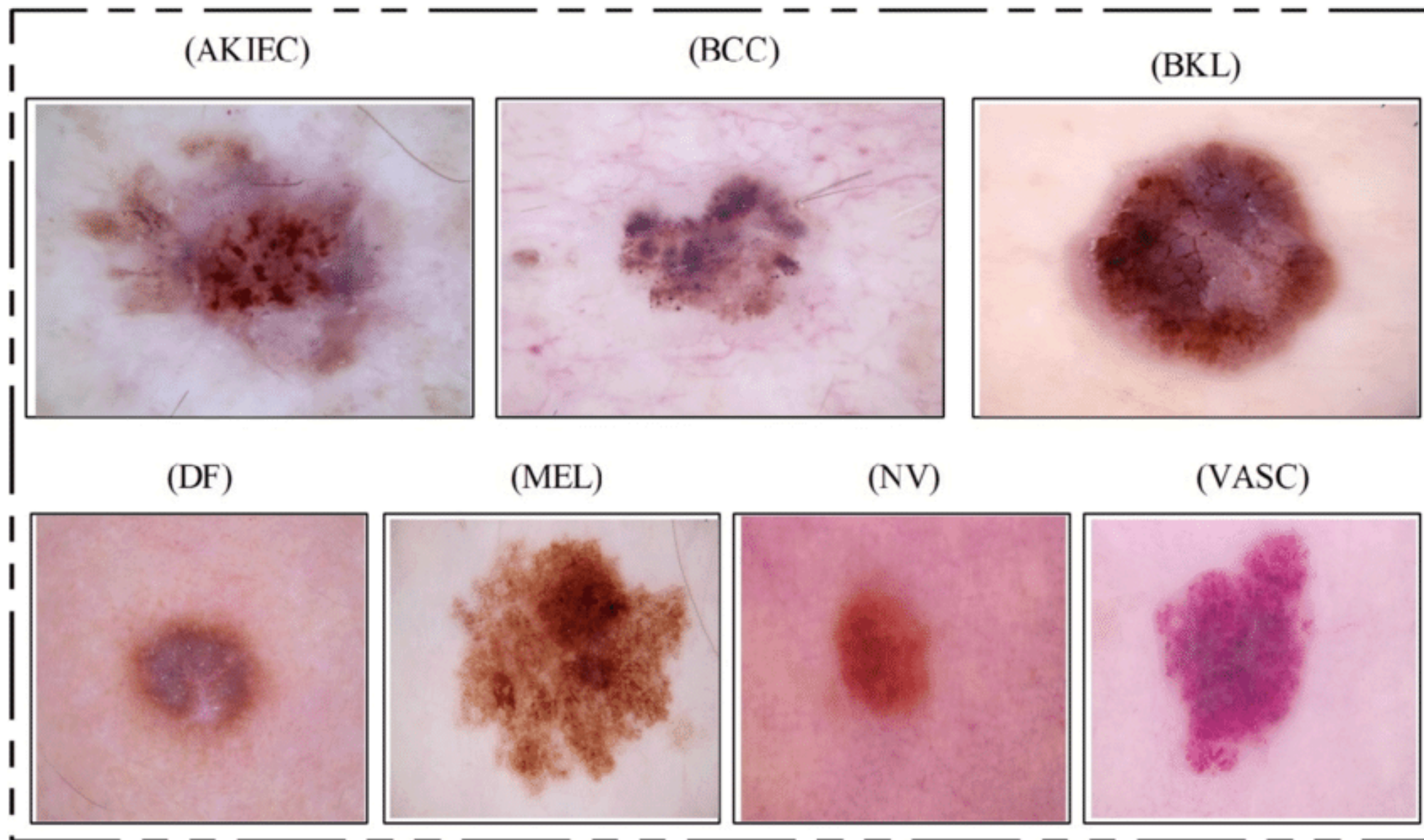
model

- ⊗ softmax and one-vs-all variants
- ⊗ both consistent w.r.t. 0-1 loss

softmax loss function

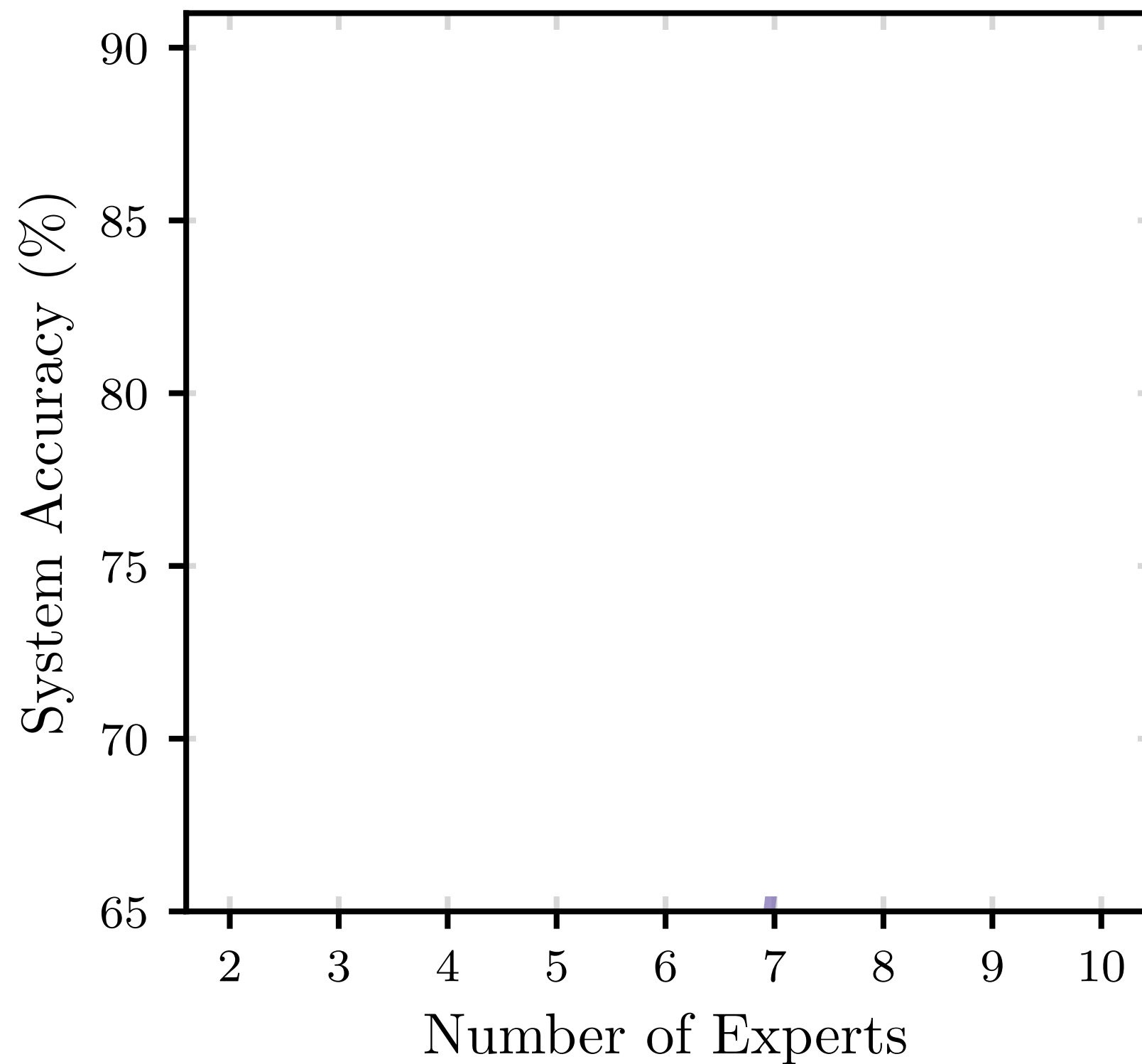
$$\ell(\theta; \mathbf{x}, y, \mathbf{m}) = -\log h_y(\mathbf{x}) - \sum_j \mathbb{I}[y = m_j] \cdot \log h_{\perp,j}(\mathbf{x})$$

# skin lesion diagnosis

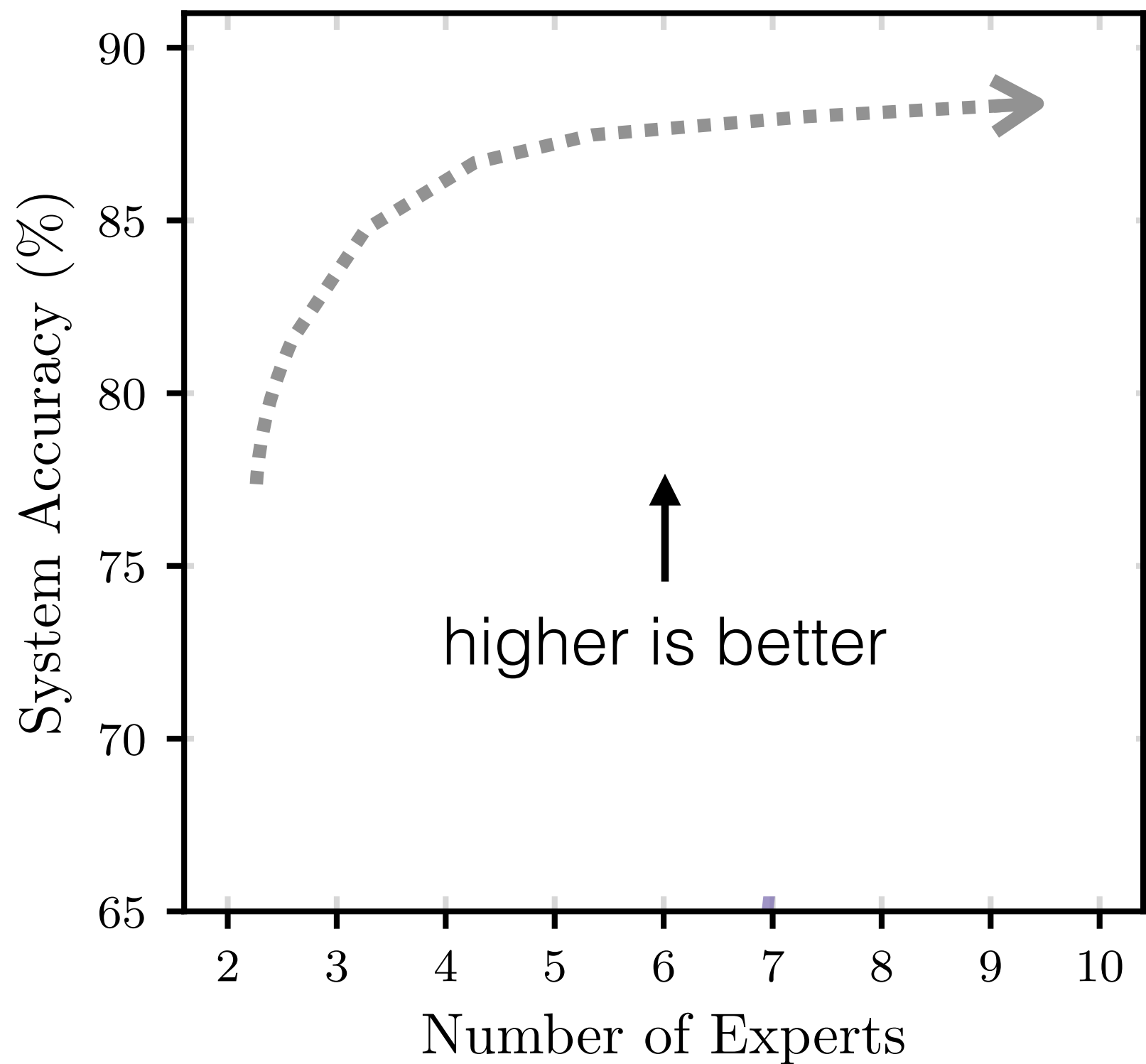




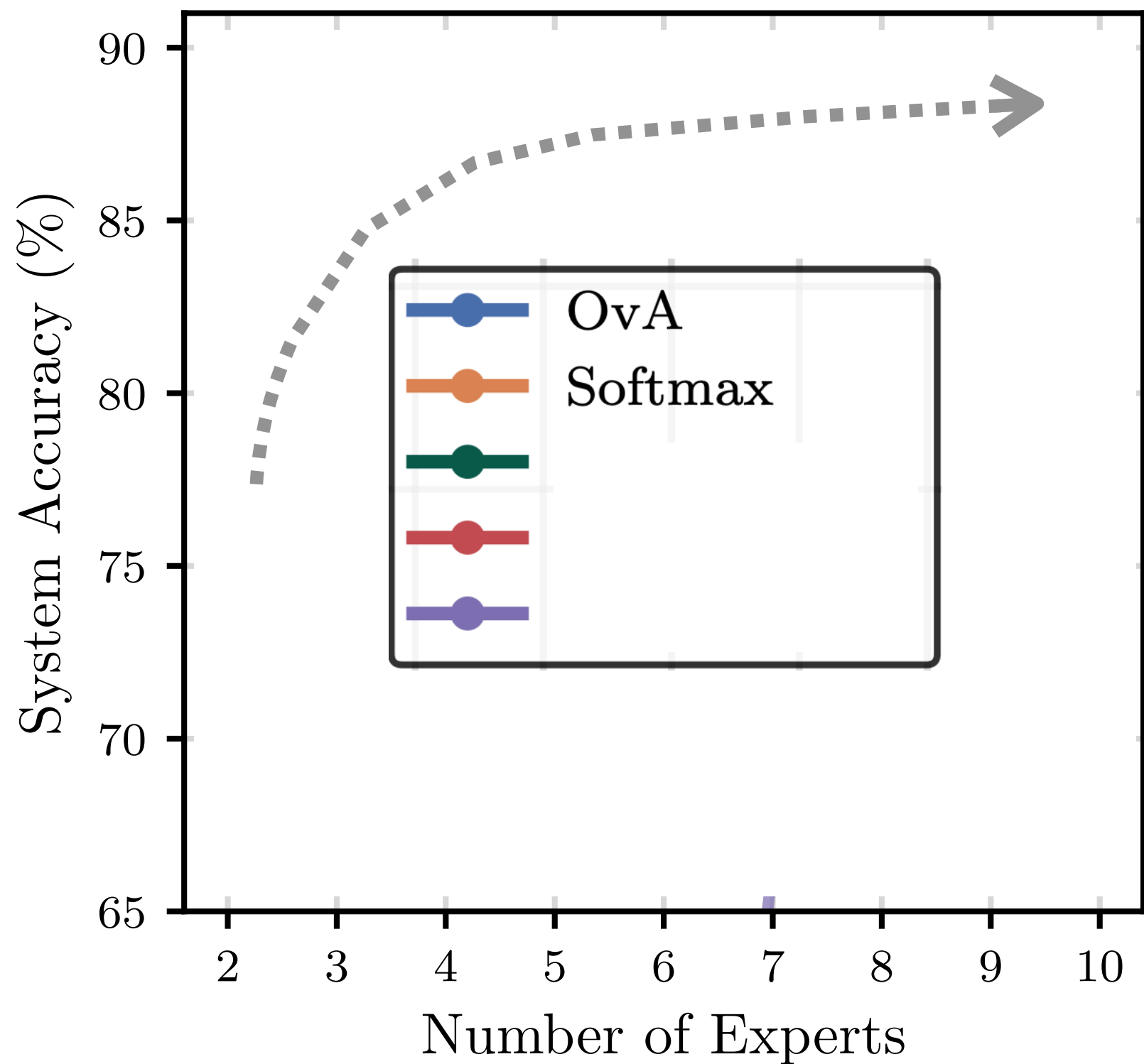
# skin lesion diagnosis



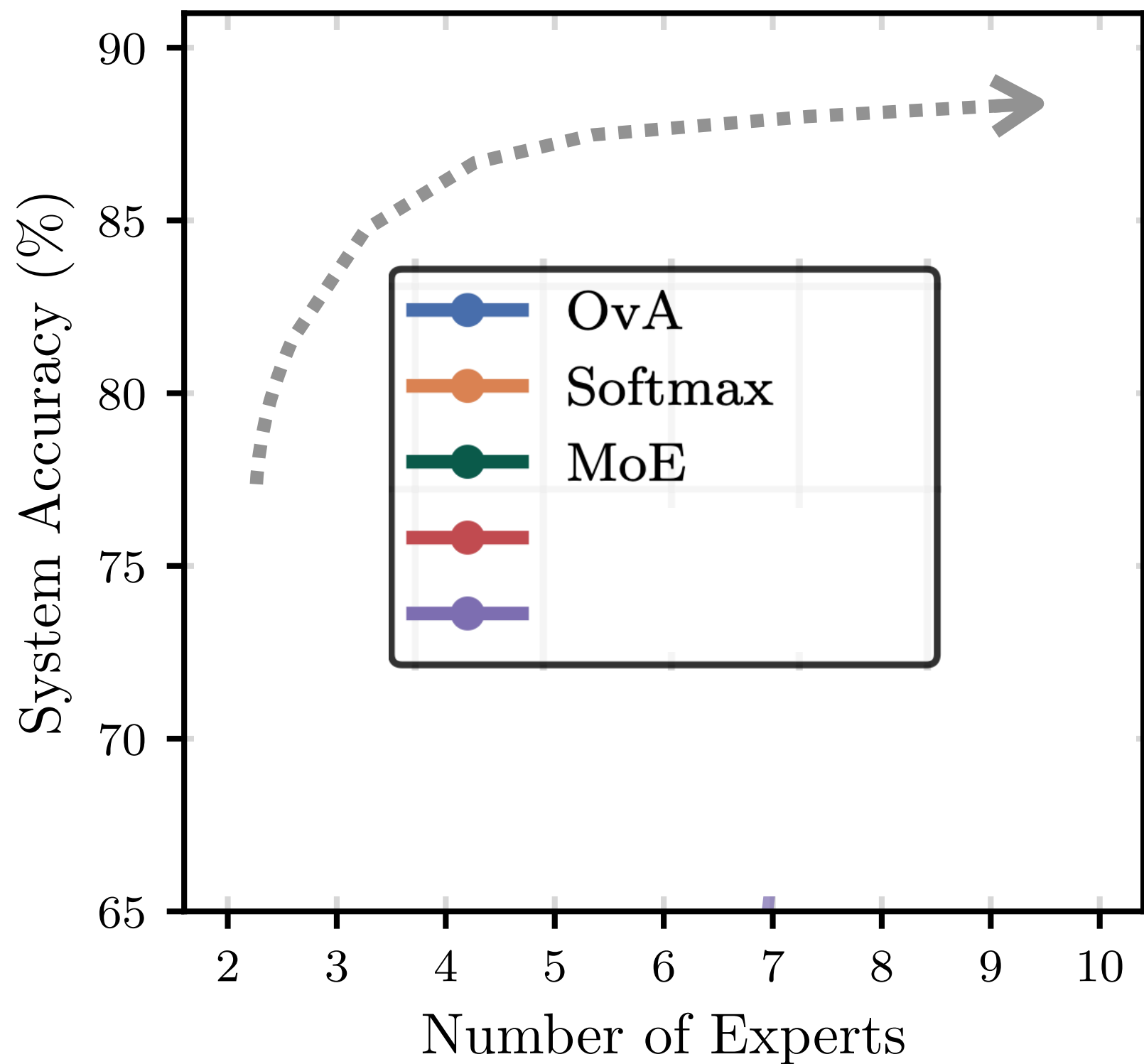
# skin lesion diagnosis



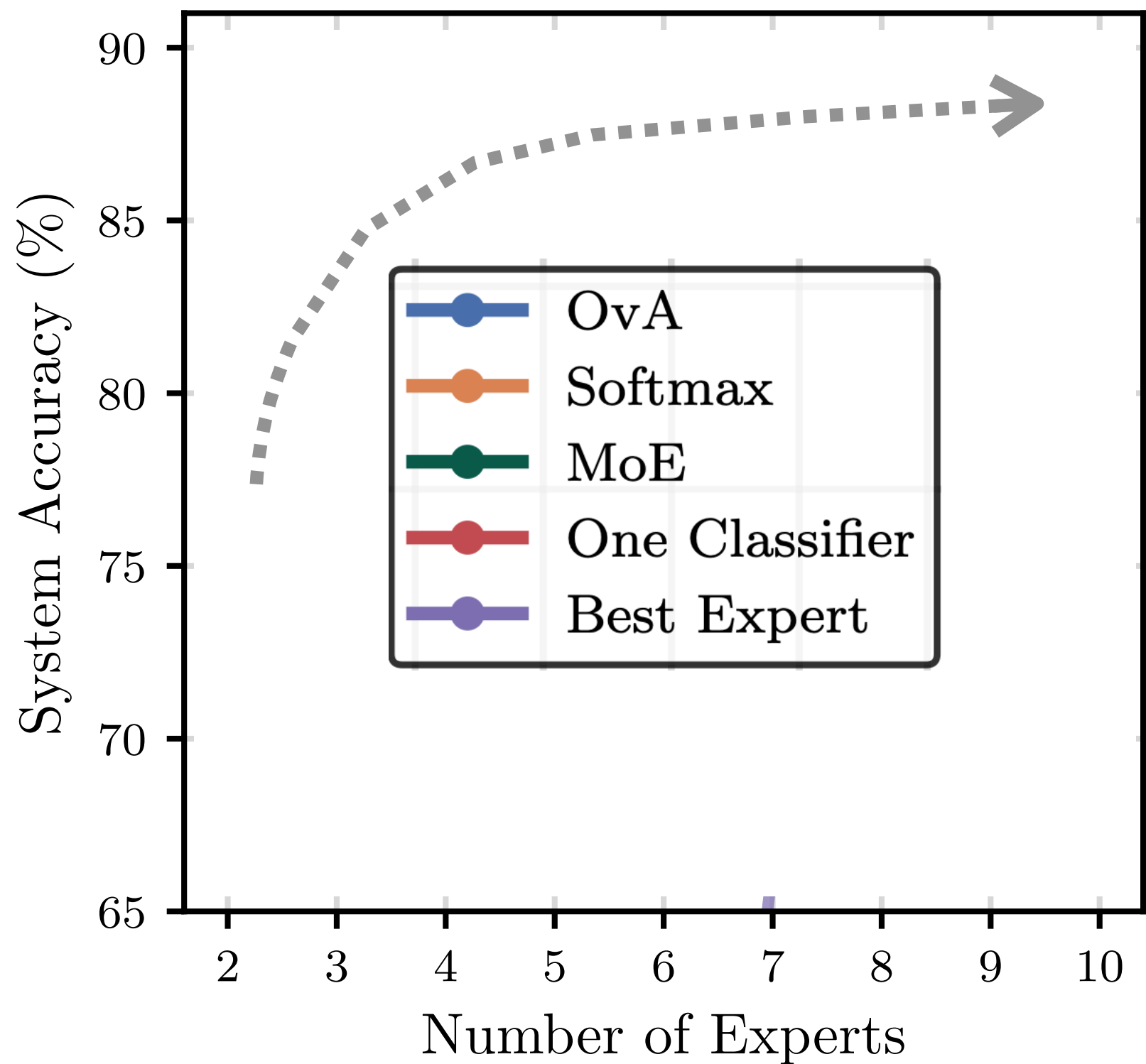
# skin lesion diagnosis



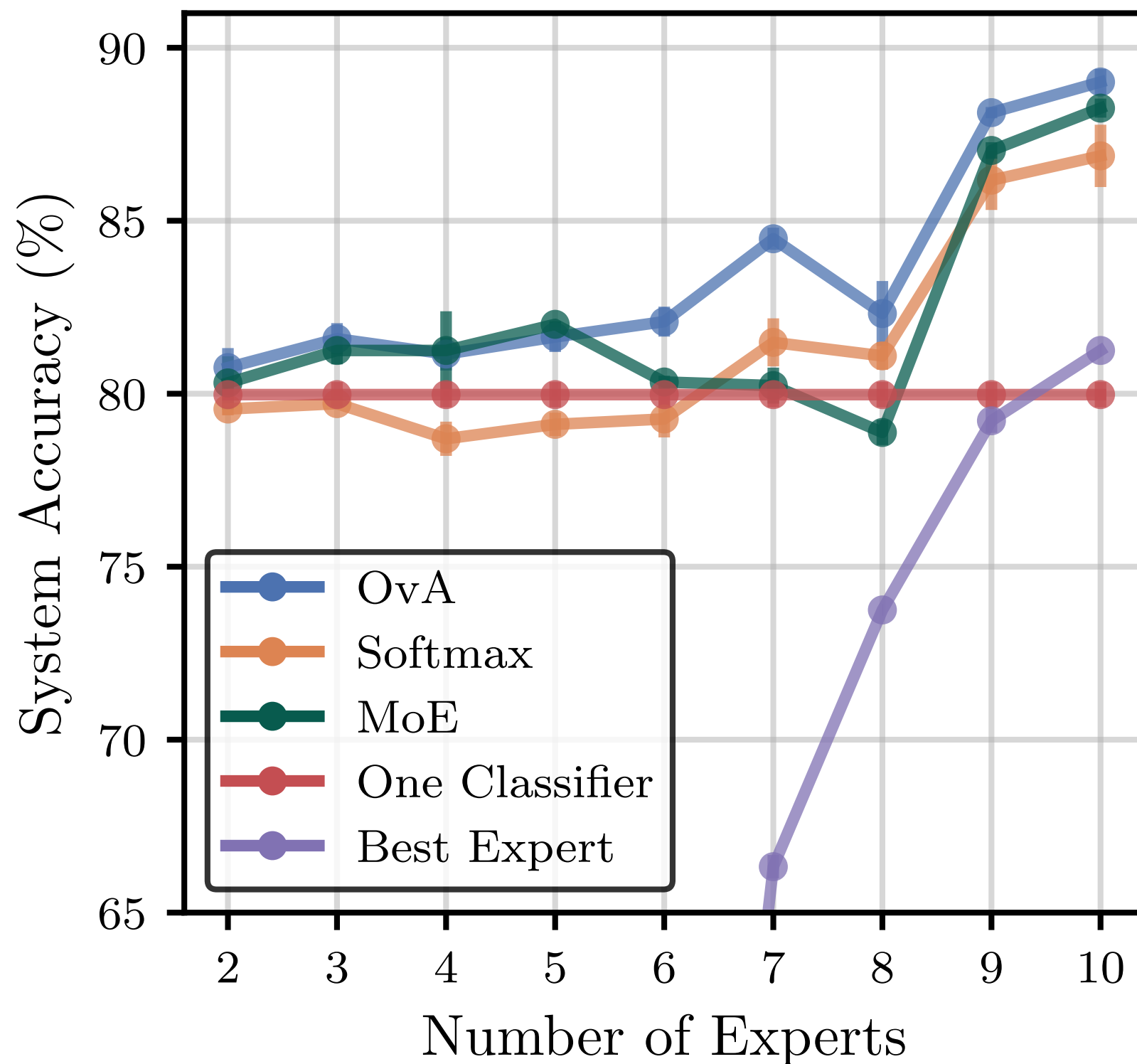
# skin lesion diagnosis



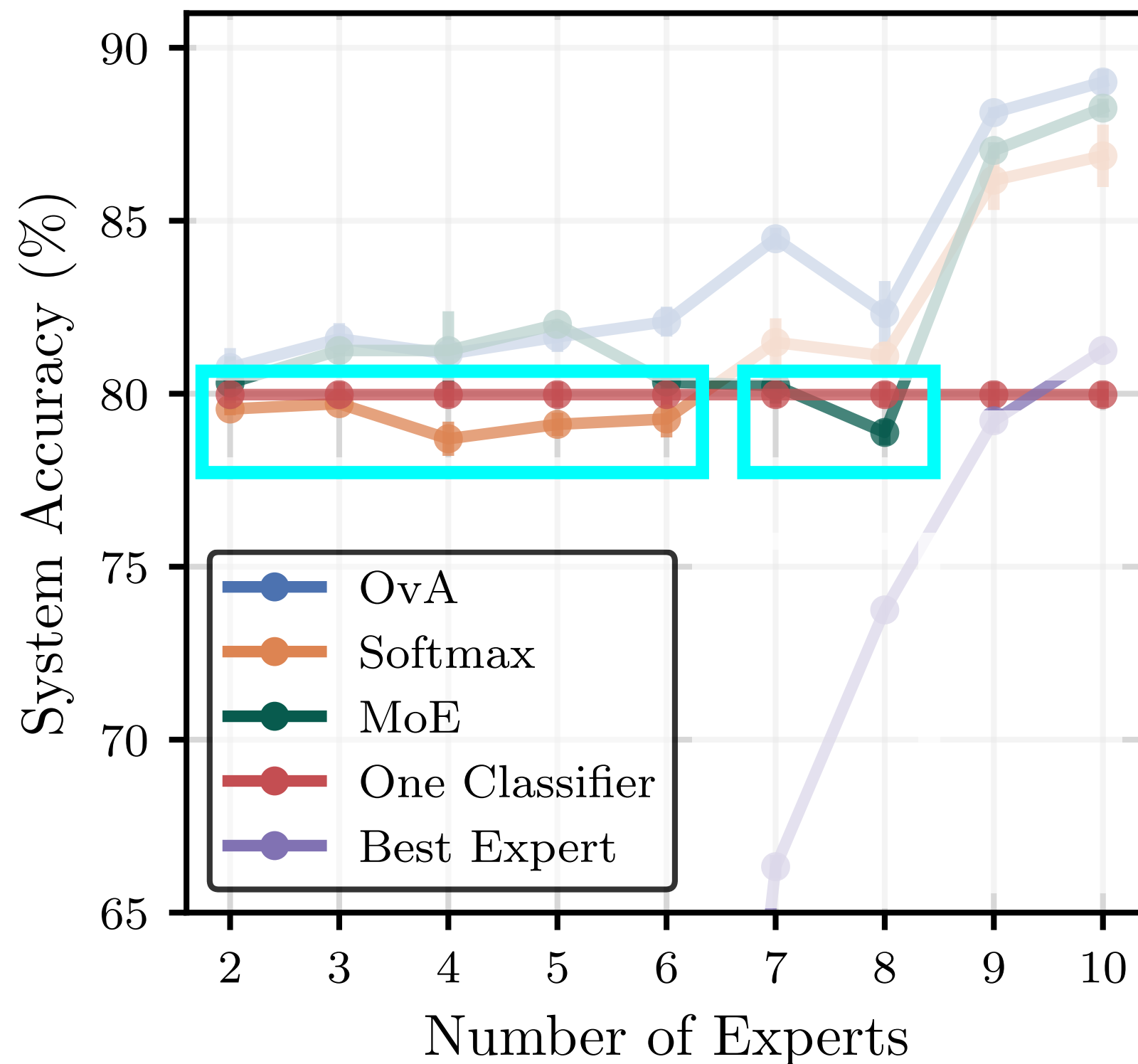
# skin lesion diagnosis



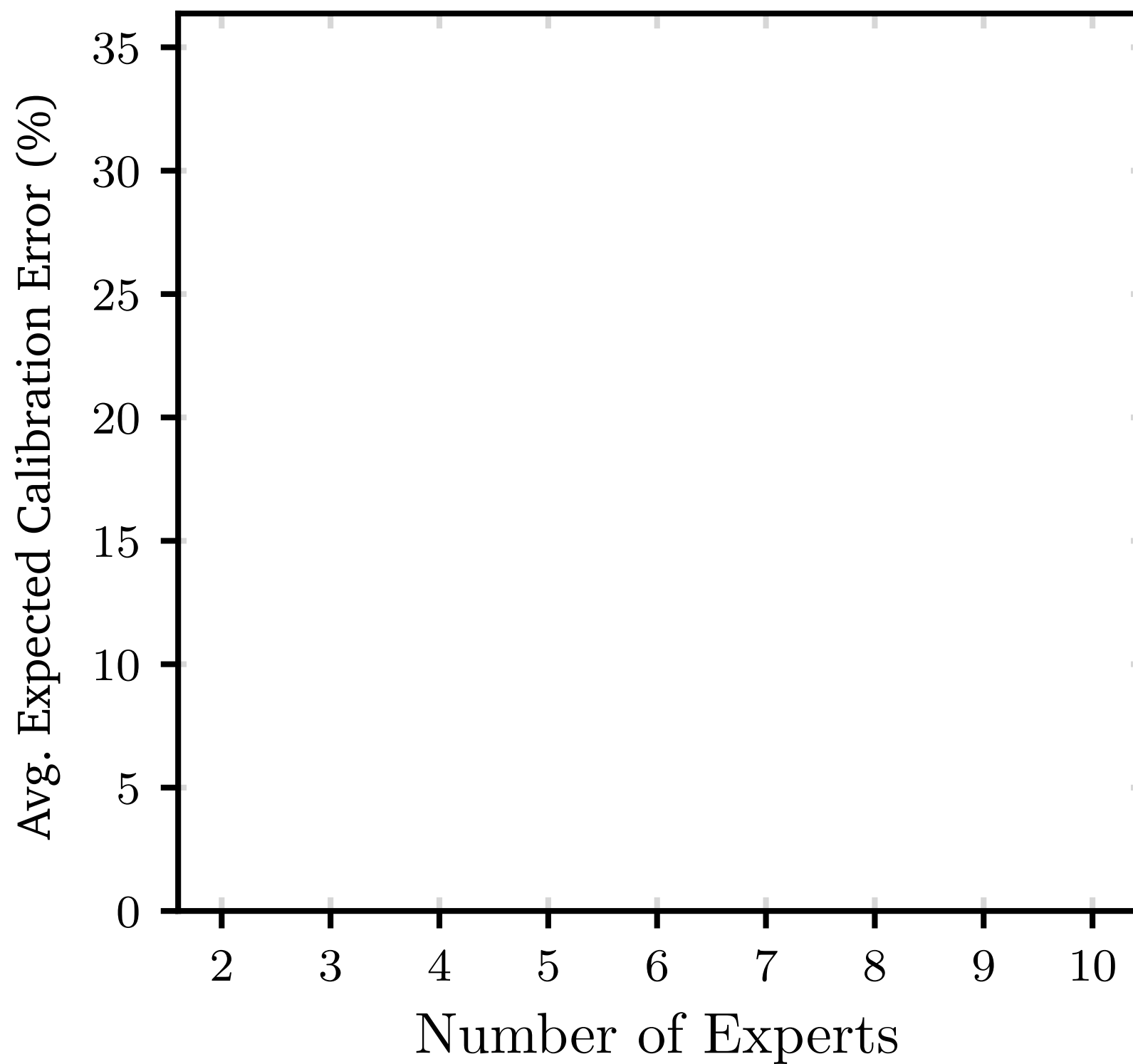
# skin lesion diagnosis



# skin lesion diagnosis

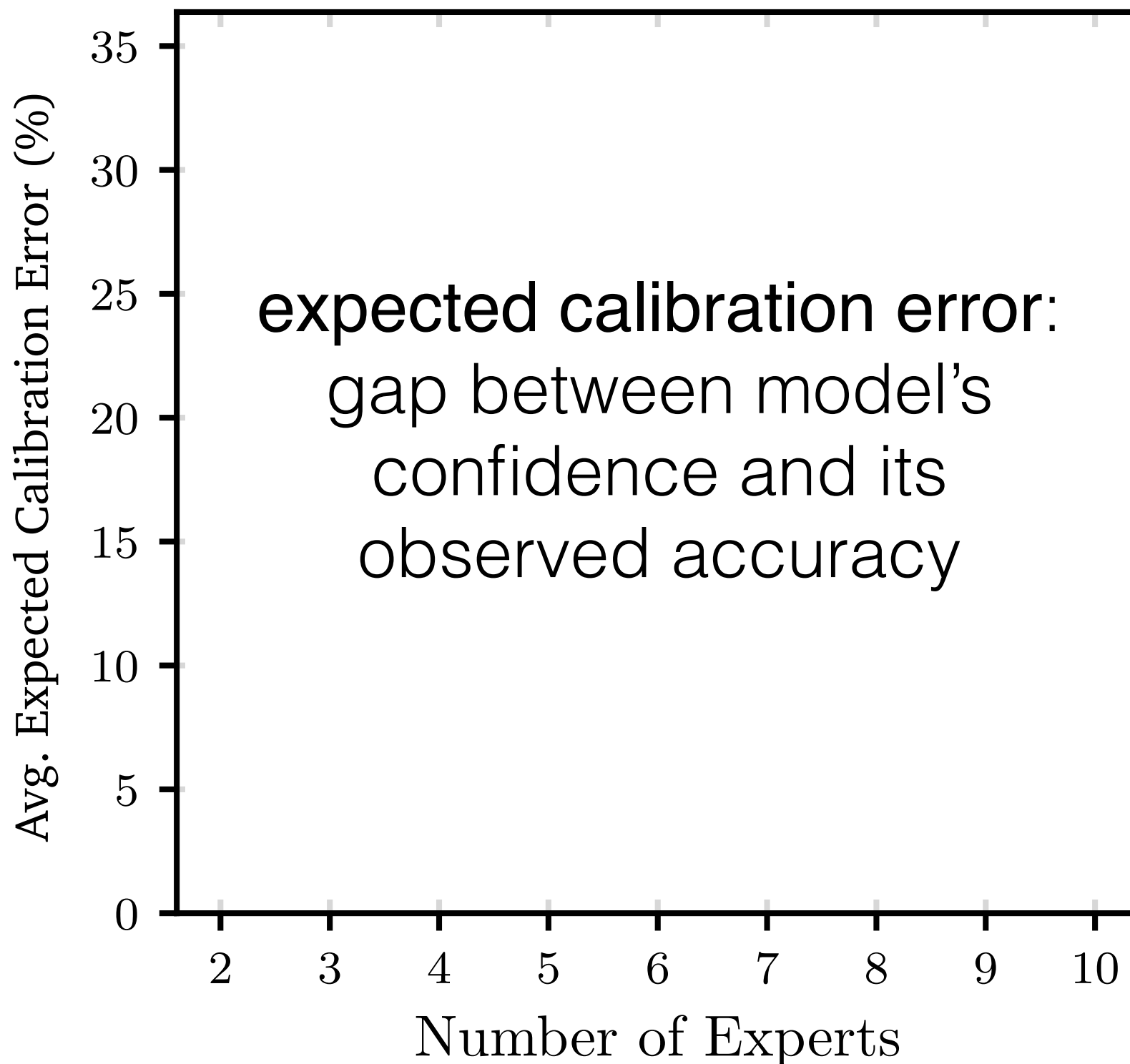


# skin lesion diagnosis

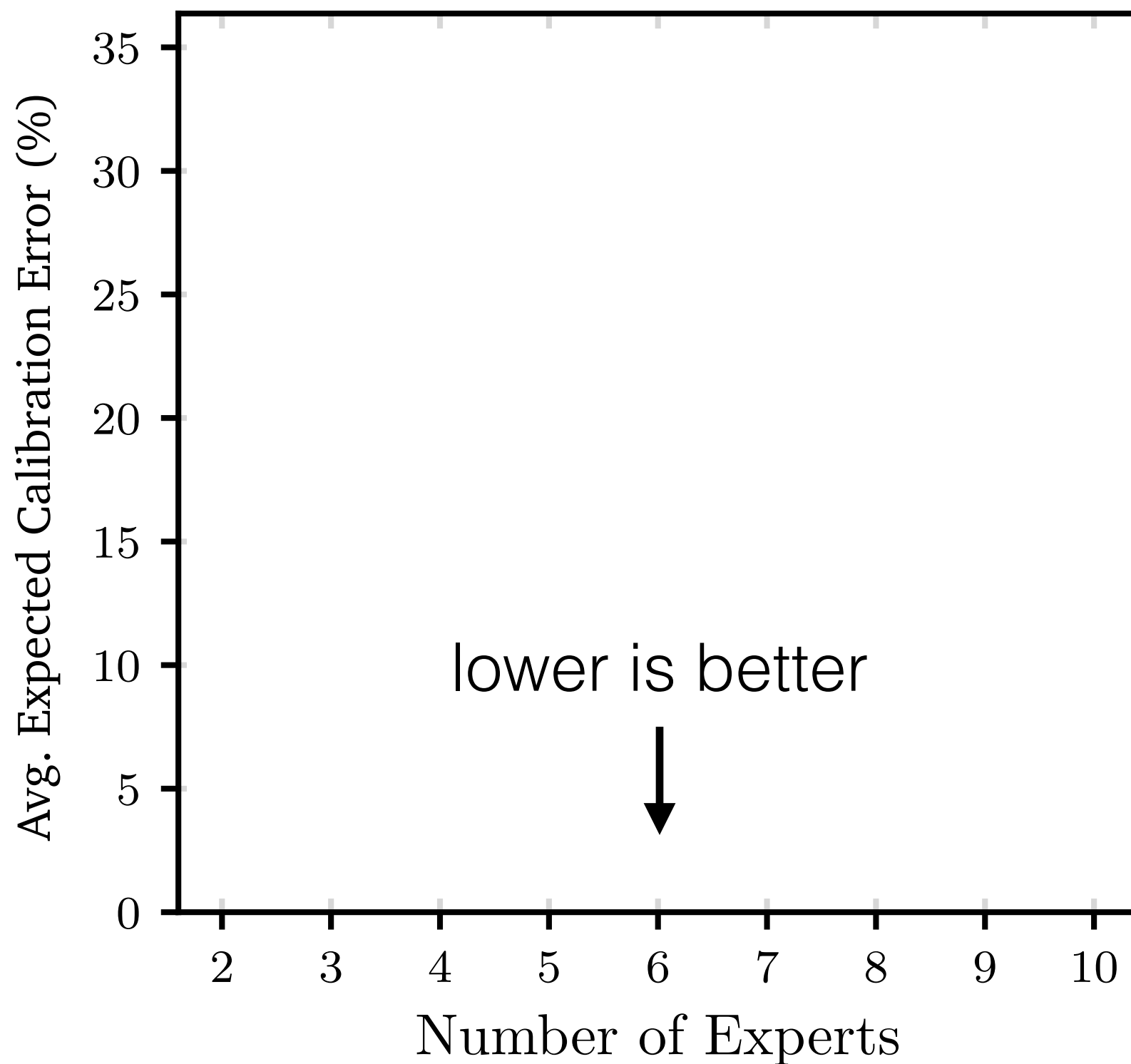




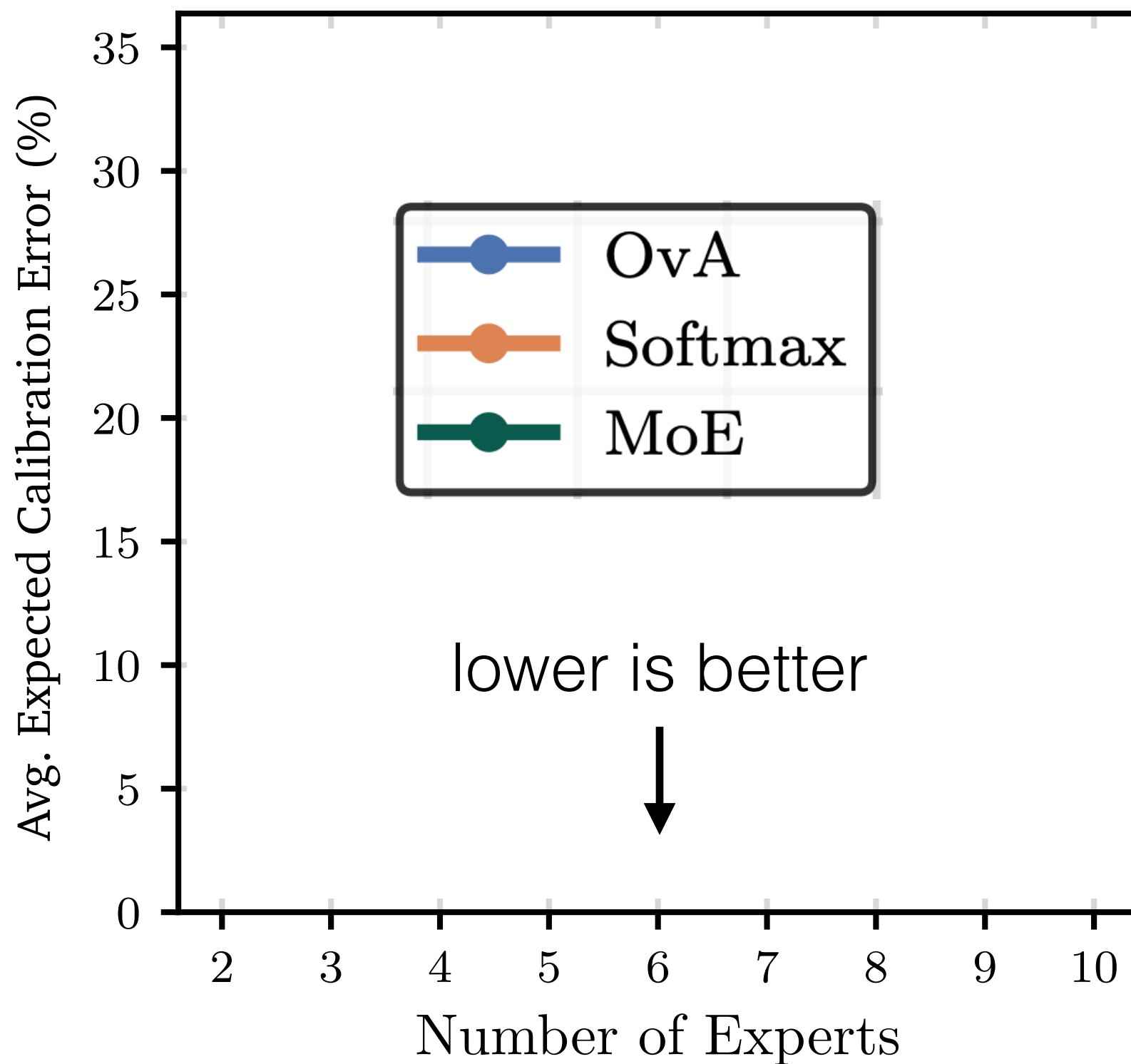
# skin lesion diagnosis



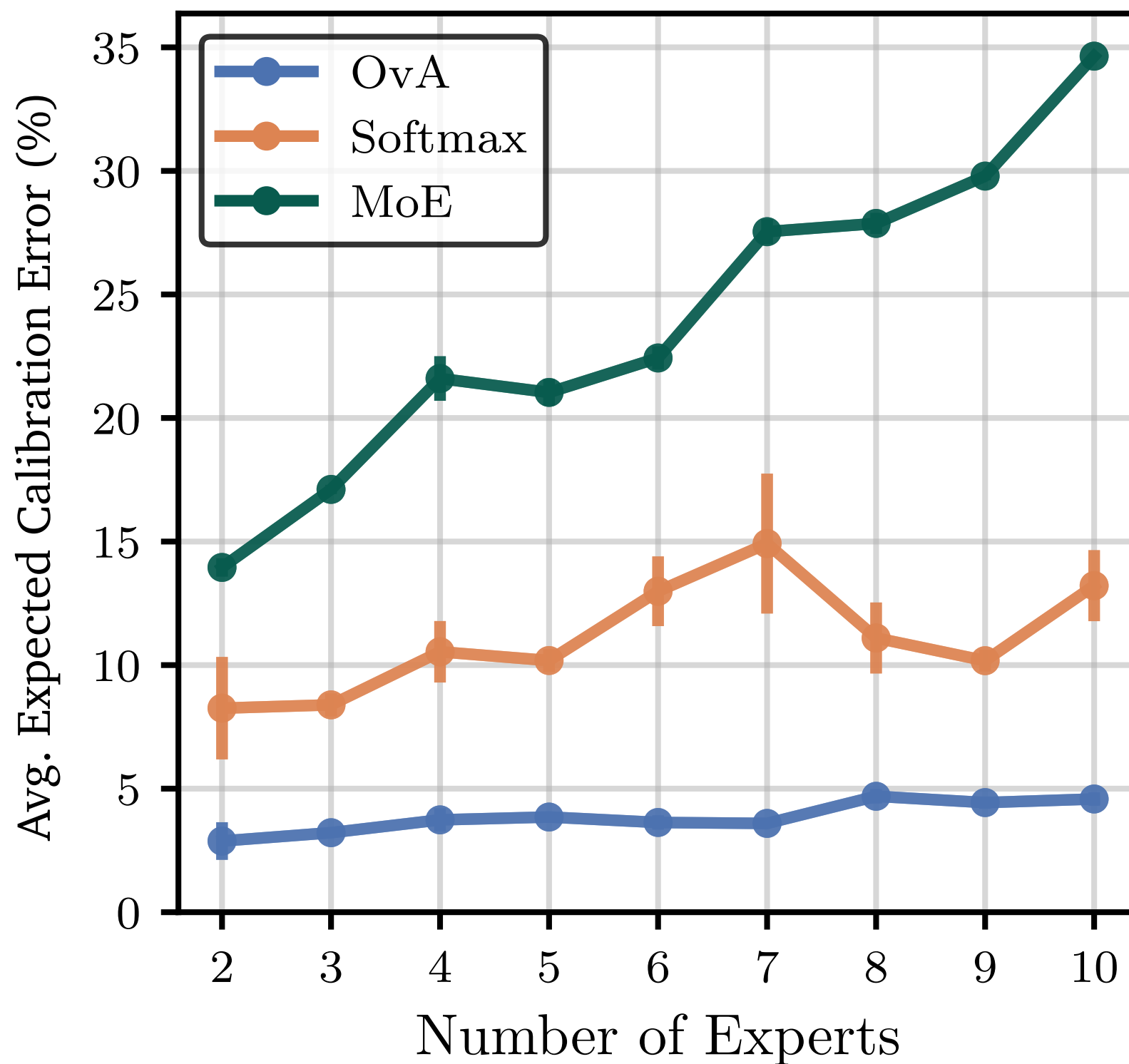
# skin lesion diagnosis



# skin lesion diagnosis



# skin lesion diagnosis



- ⊗ single expert
  - ⊗ softmax surrogate loss
  - ⊗ improving calibration via one-vs-all

- ⊗ **multiple experts**

- ⊗ surrogate losses
- ⊗ conformal sets of experts

- ⊗ population of experts
  - ⊗ surrogate losses
  - ⊗ meta-learning a rejector

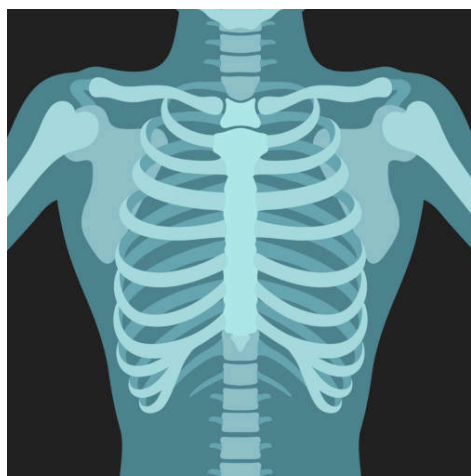
- ⊗ single expert
  - ⊗ softmax surrogate loss
  - ⊗ improving calibration via one-vs-all

- ⊗ **multiple experts**
  - ⊗ surrogate losses

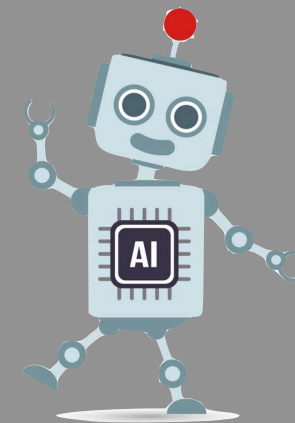
- ⊗ conformal sets of experts

- ⊗ population of experts
  - ⊗ surrogate losses
  - ⊗ meta-learning a rejector

input  
features



allocation  
mechanism



classifier



expert #1

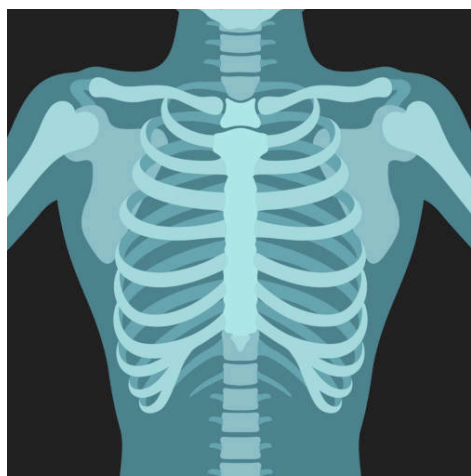


expert #3

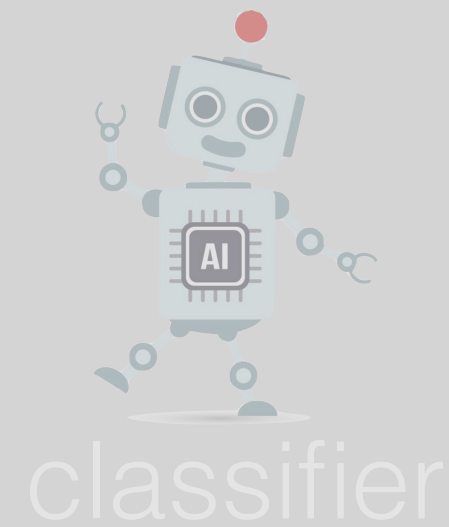


expert #2

input  
features



allocation  
mechanism



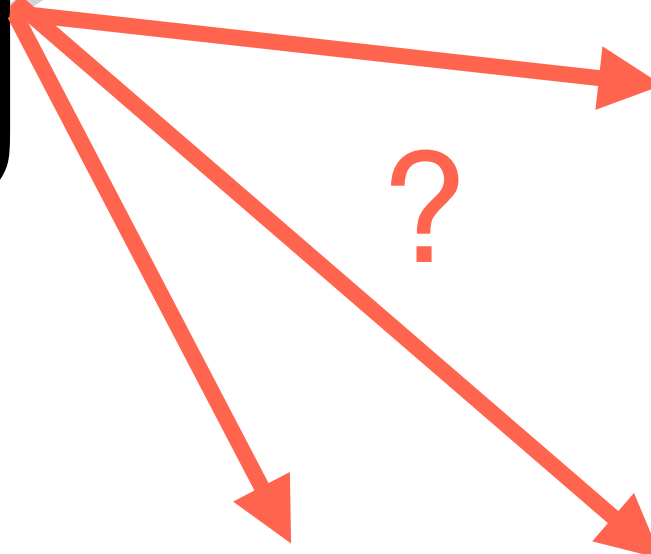
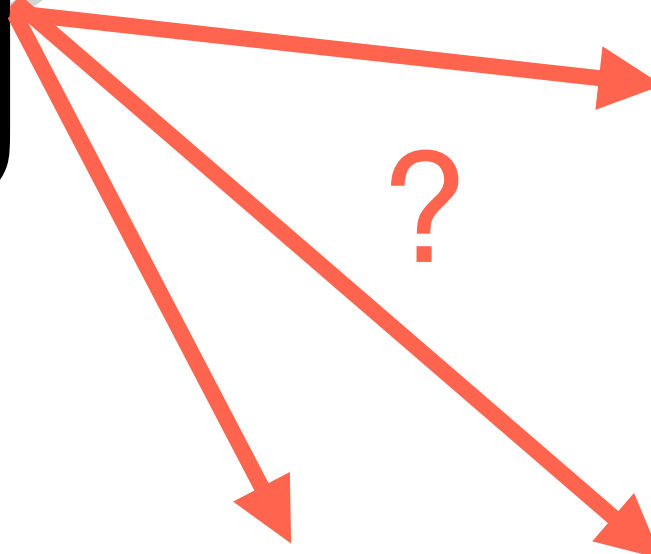
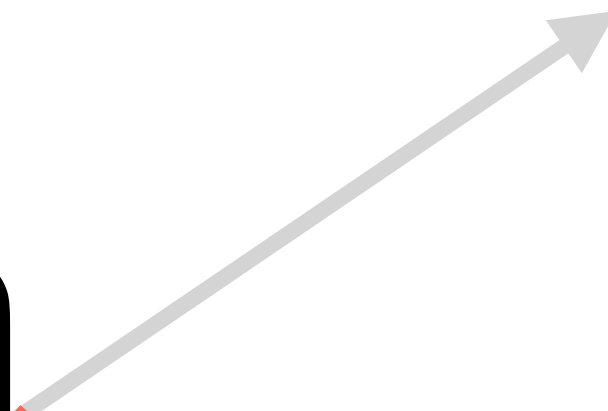
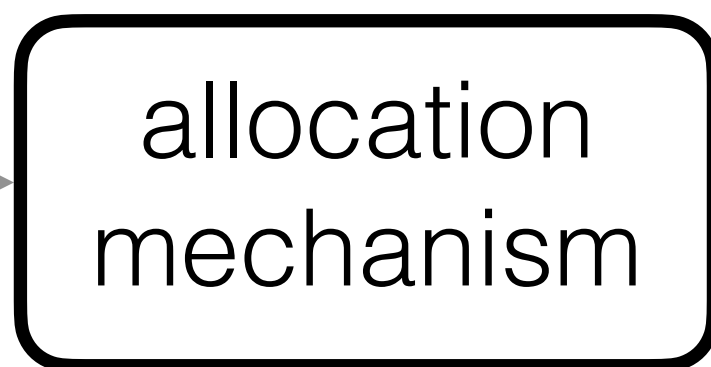
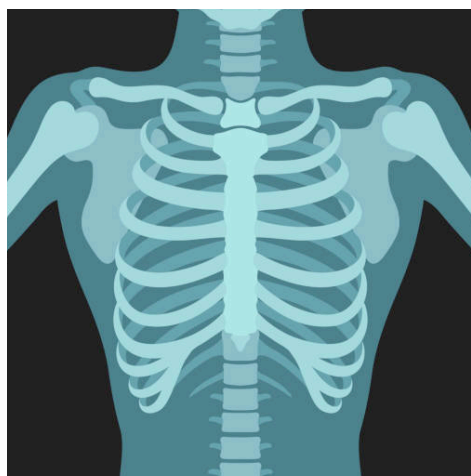
expert #1

expert #3

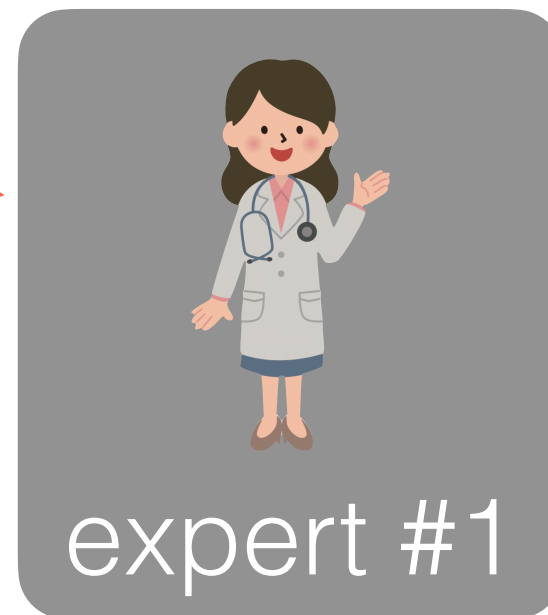
expert #2



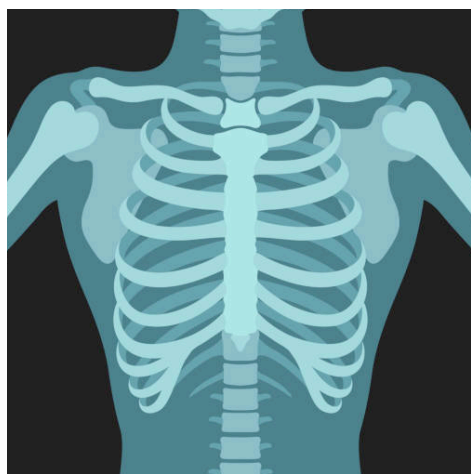
input  
features



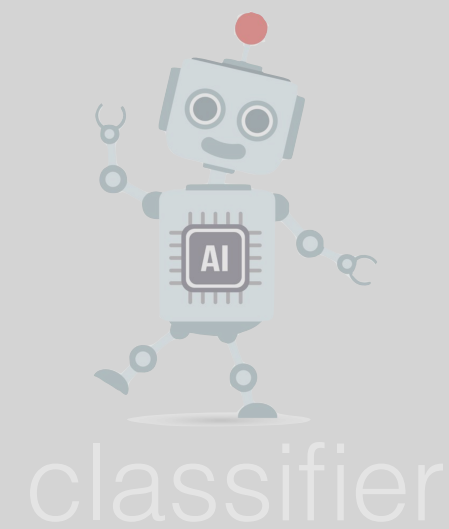
?



input  
features



allocation  
mechanism



conformal inference

# conformal inference

assume there's a best expert,  $j^*$ :

$$\mathbb{P}(m_{j^*} = y | \mathbf{x}) > \mathbb{P}(m_e = y | \mathbf{x}), \quad \forall e \neq j^*$$

# conformal inference

assume there's a best expert,  $j^*$ :

$$\mathbb{P}(m_{j^*} = y | \mathbf{x}) > \mathbb{P}(m_e = y | \mathbf{x}), \quad \forall e \neq j^*$$

construct a confidence set of experts:

$$\mathbb{P}(j^* \in C(\mathbf{x})) \geq 1 - \alpha$$

# conformal inference

assume there's a best expert,  $j^*$ :

$$\mathbb{P}(m_{j^*} = y | \mathbf{x}) > \mathbb{P}(m_e = y | \mathbf{x}), \quad \forall e \neq j^*$$

construct a confidence set of experts:

$$\mathbb{P}(j^* \in C(\mathbf{x})) \geq 1 - \alpha$$

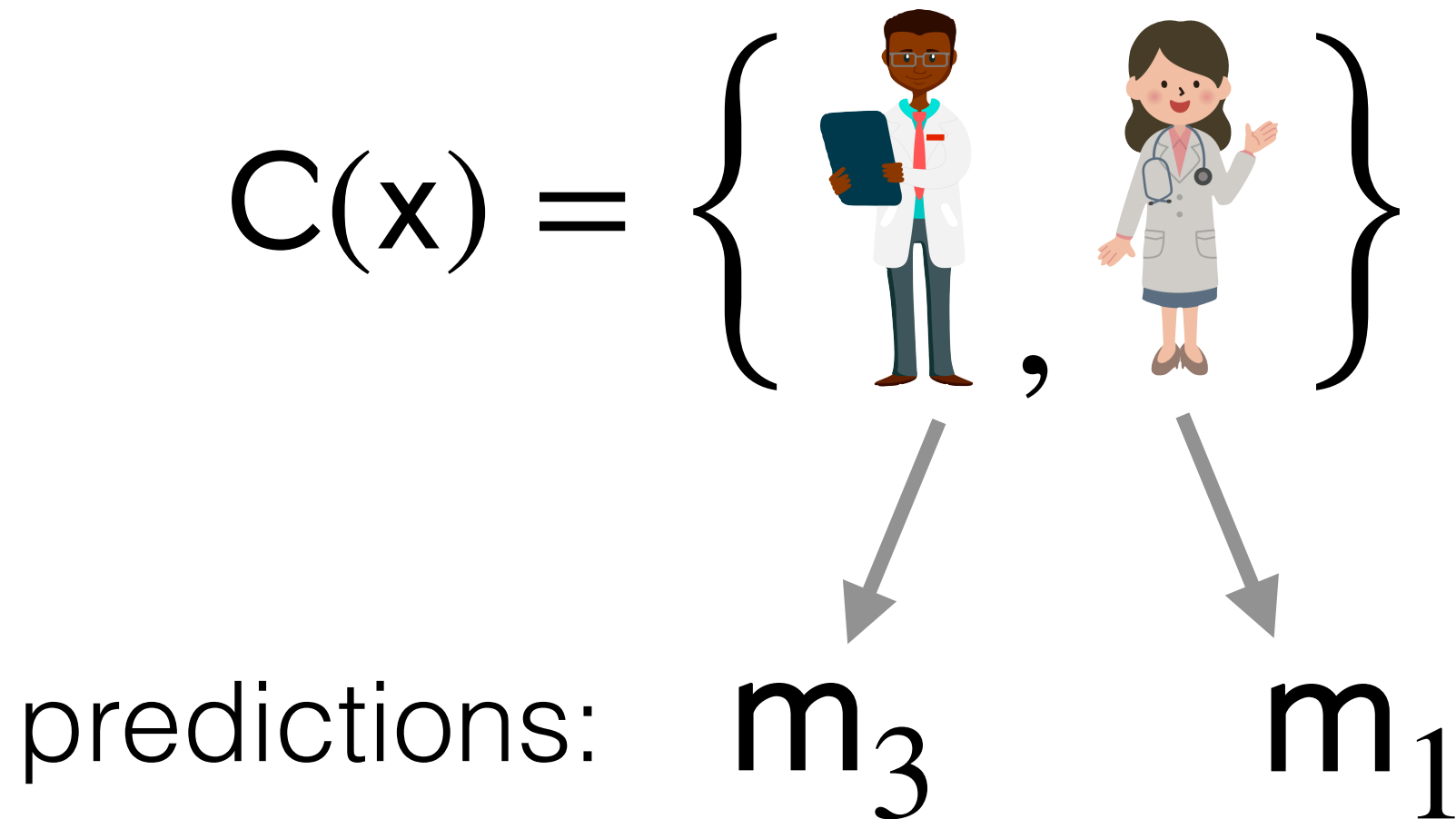


*team of experts: adaptive in size and membership*

# conformal inference: ensembling

$$C(x) = \left\{ \text{doctor}, \text{nurse} \right\}$$

# conformal inference: ensembling





# conformal inference: ensembling

$$C(x) = \left\{ \text{doctor}_3, \text{doctor}_1 \right\}$$

predictions:

$m_3$

$m_1$

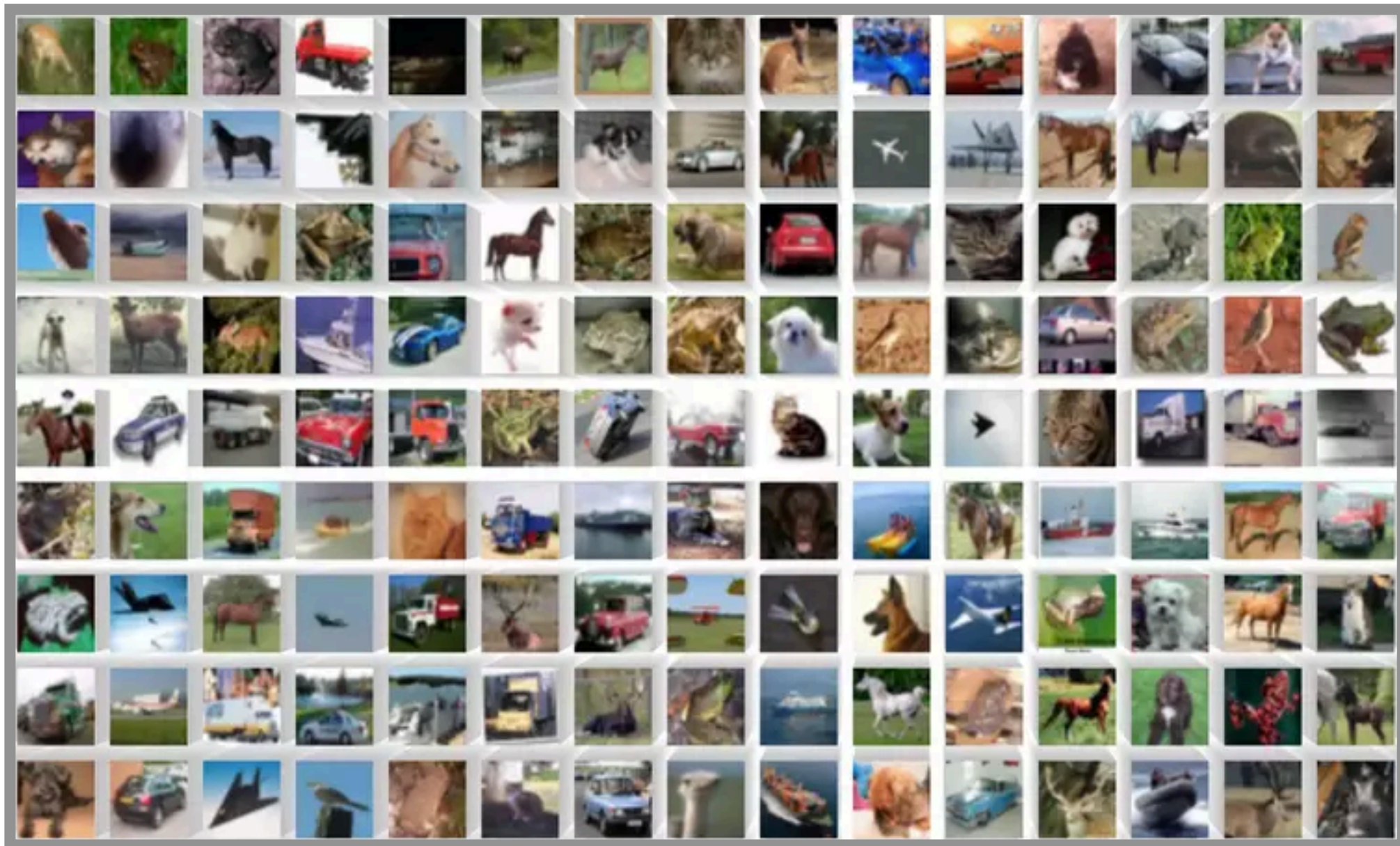
aggregated prediction:

$\bar{m}$

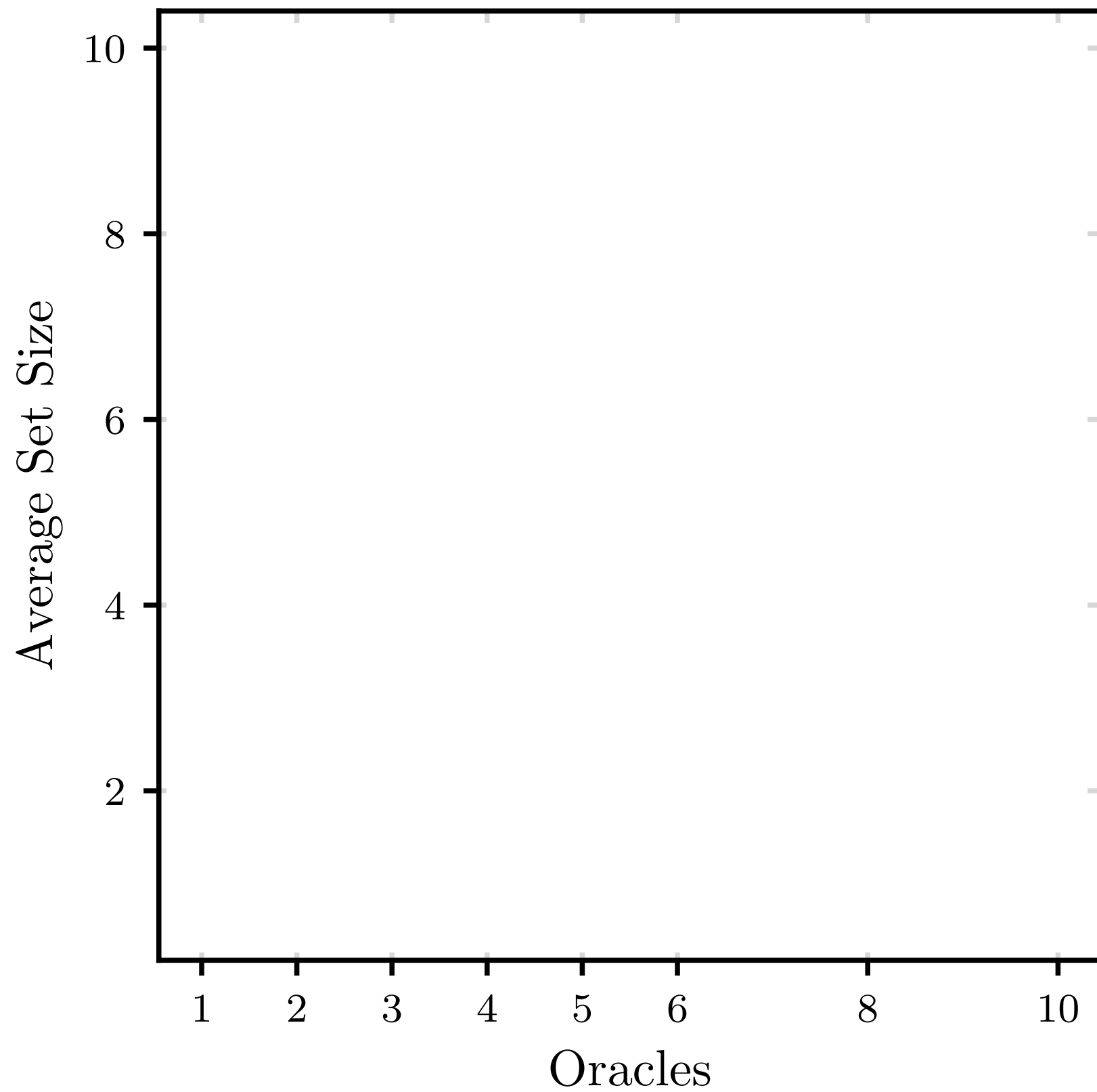
(e.g. by voting)

# expert selection

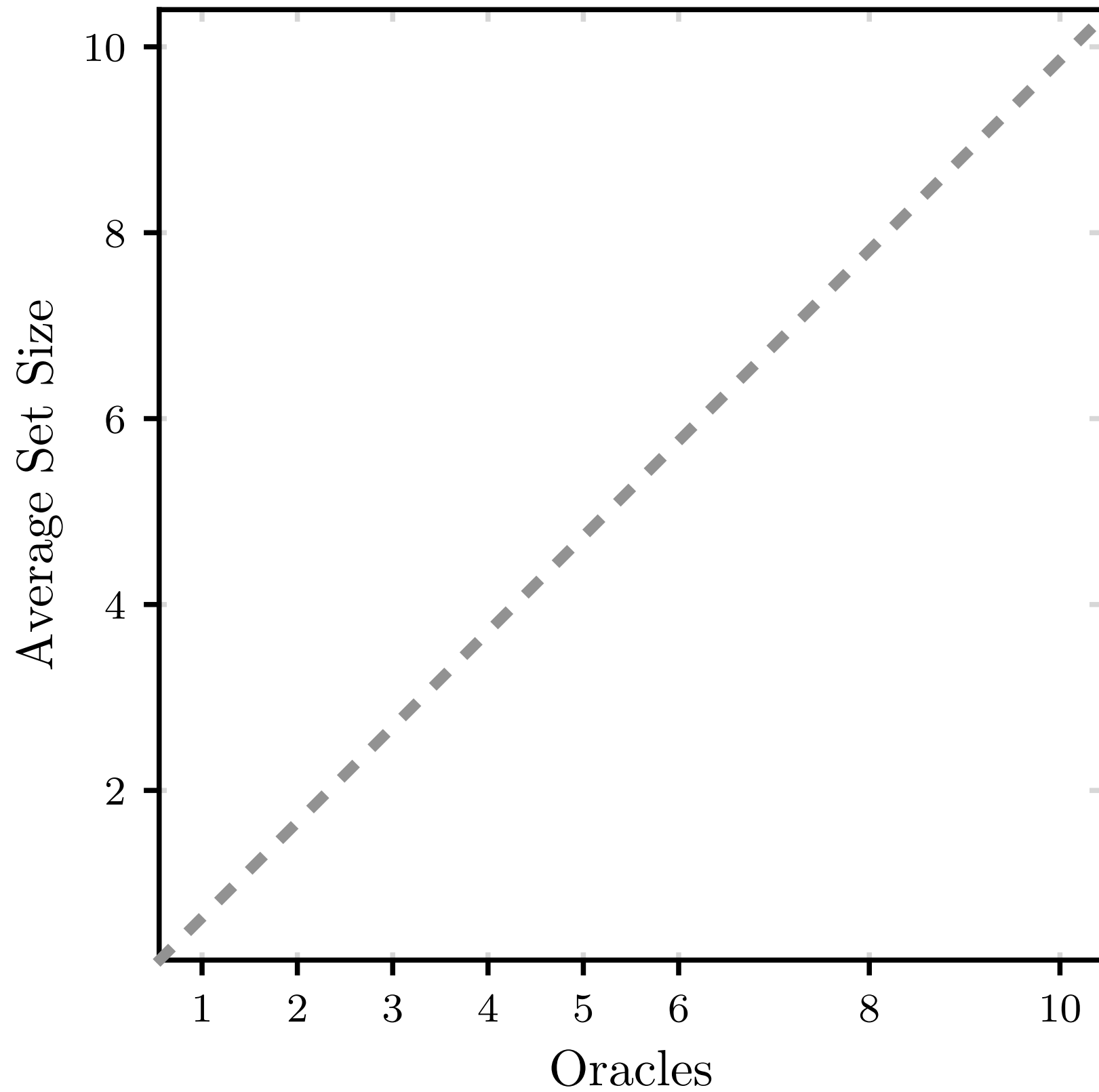
## CIFAR-10



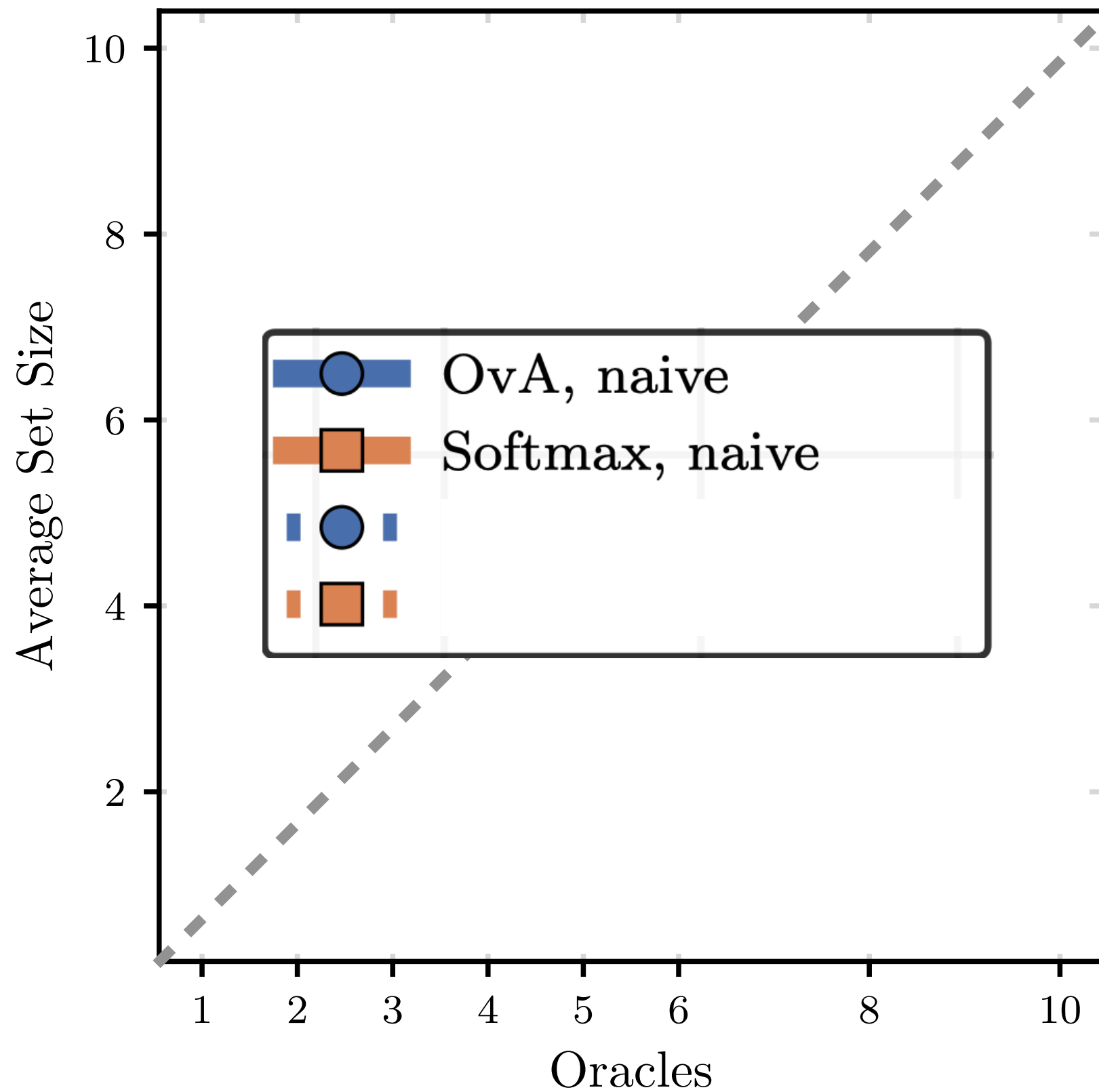
# expert selection



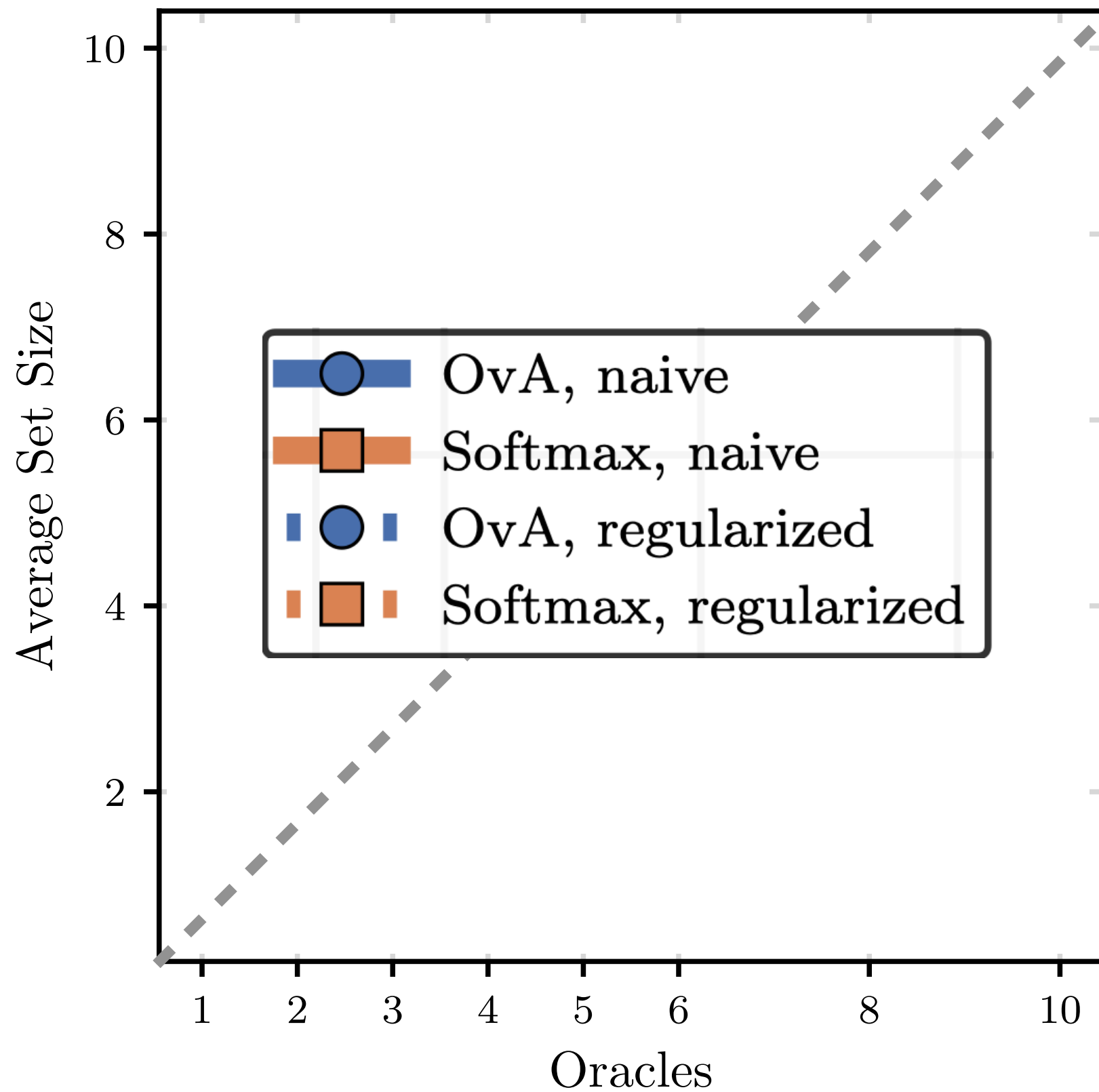
# expert selection



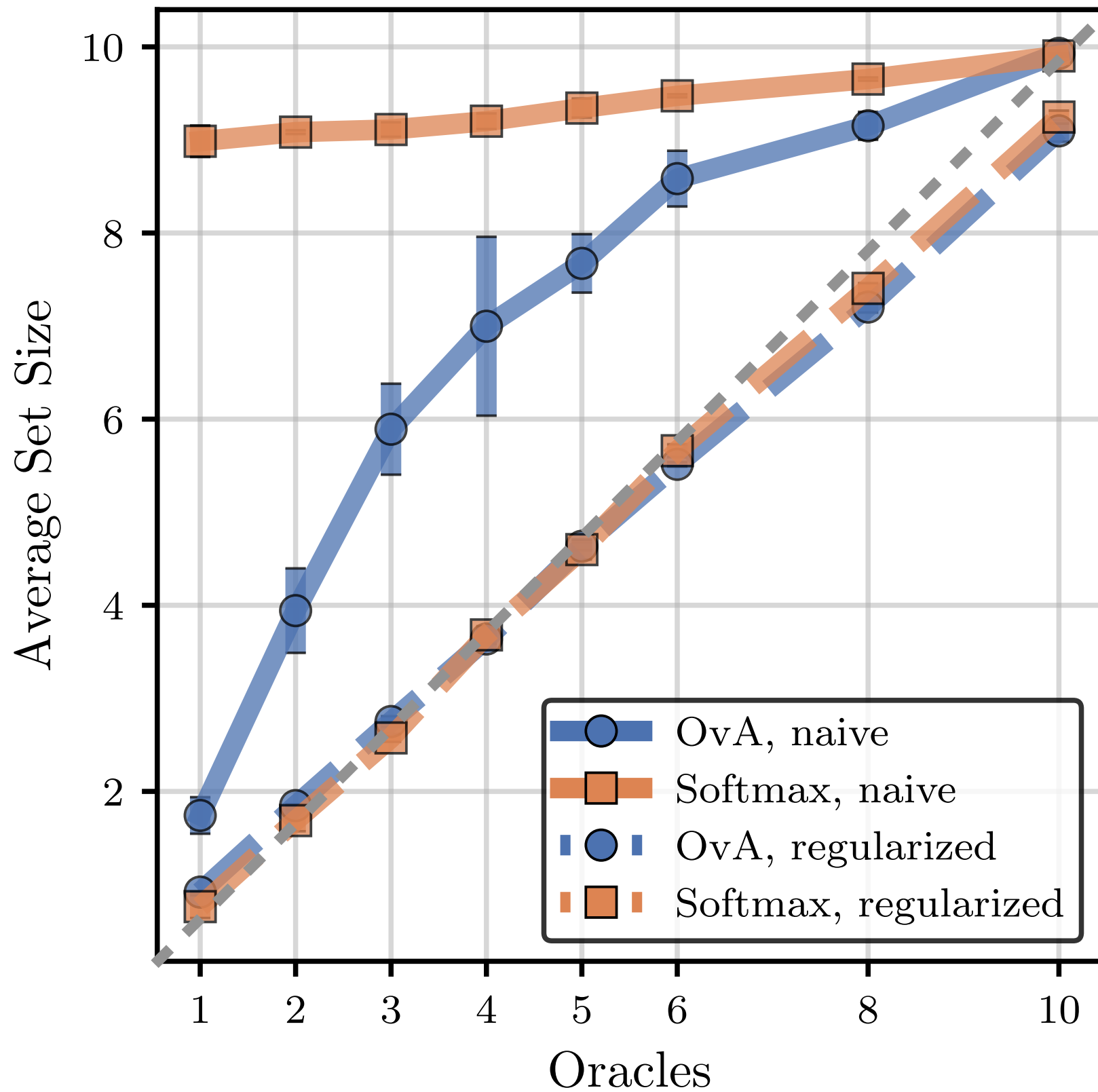
# expert selection



# expert selection



# expert selection



- ⊗ single expert
  - ⊗ softmax surrogate loss
  - ⊗ improving calibration via one-vs-all

- ⊗ **multiple experts**
  - ⊗ surrogate losses
  - ⊗ conformal sets of experts

- ⊗ population of experts
  - ⊗ surrogate losses
  - ⊗ meta-learning a rejector

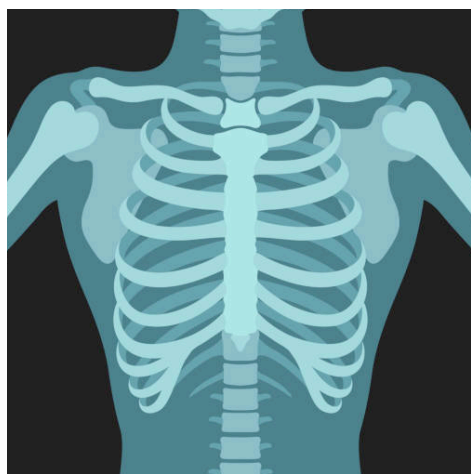


- ⊗ single expert
  - ⊗ softmax surrogate loss
  - ⊗ improving calibration via one-vs-all

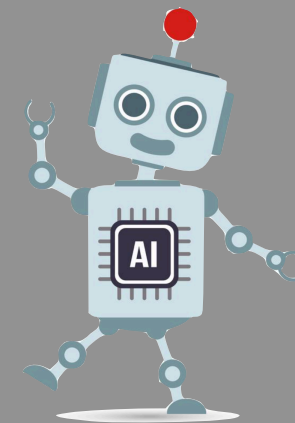
- ⊗ multiple experts
  - ⊗ surrogate losses
  - ⊗ conformal sets of experts

- ⊗ **population of experts**
  - ⊗ surrogate losses
  - ⊗ meta-learning a rejector

input  
features



allocation  
mechanism



classifier



expert #1

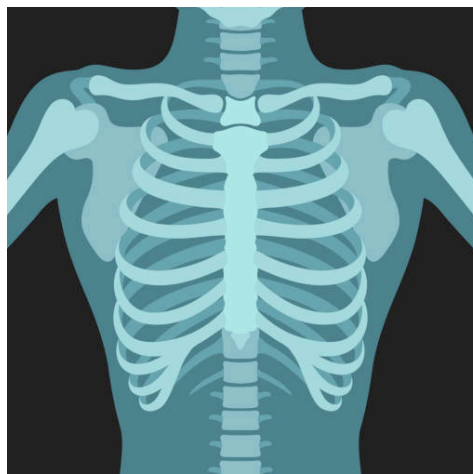


expert #3



expert #2

input  
features



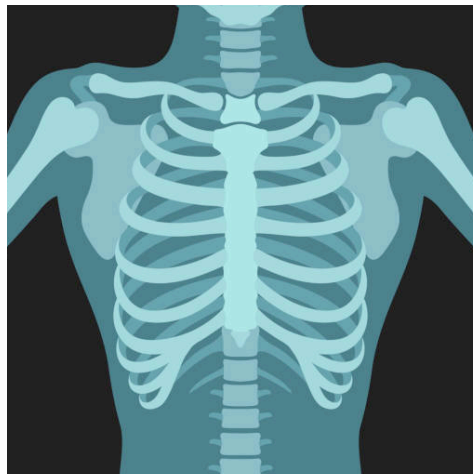
allocation  
mechanism



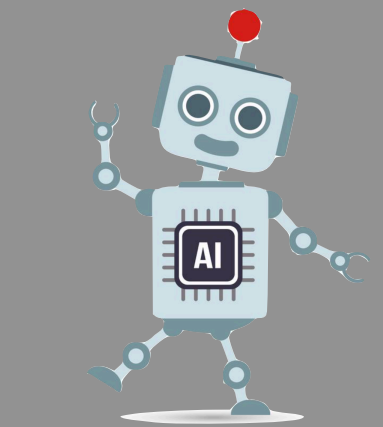
limitations?



input  
features



allocation  
mechanism

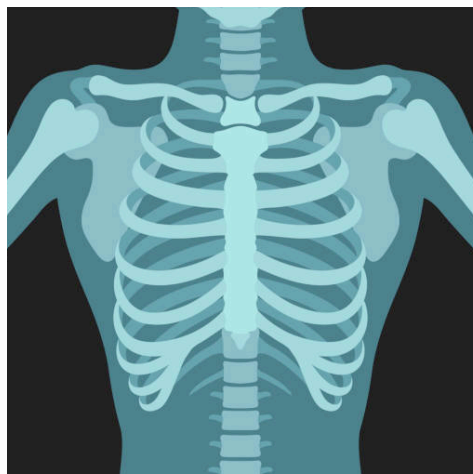


classifier

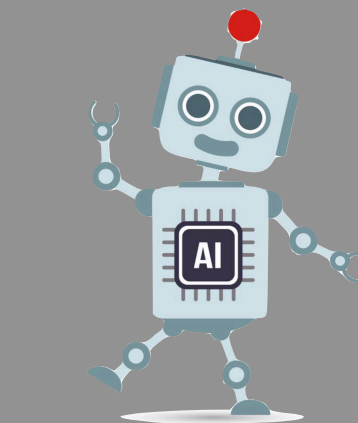
?

expert

input  
features

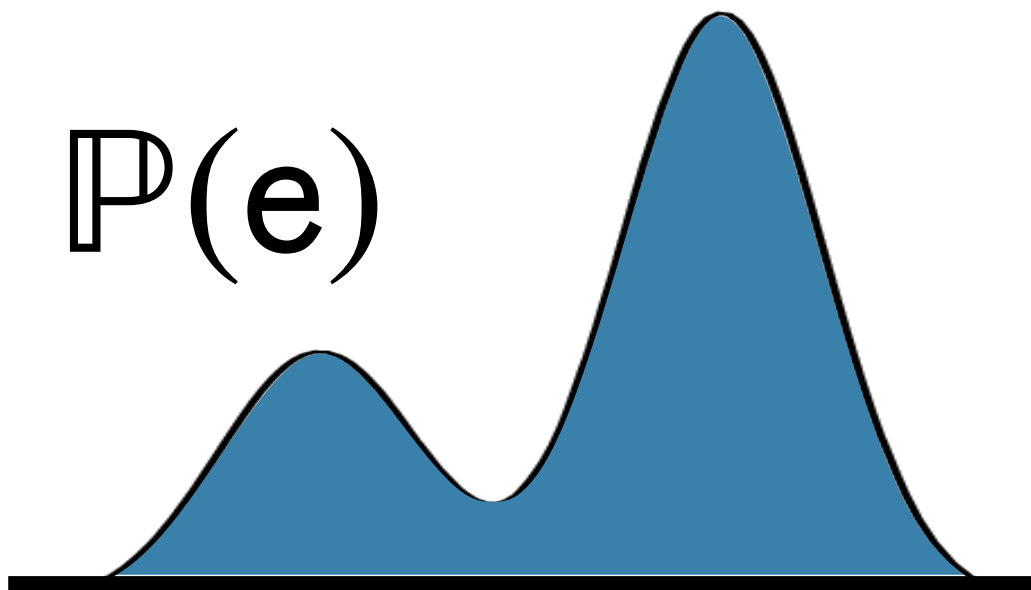


allocation  
mechanism



classifier

$\mathbb{P}(\mathbf{e})$

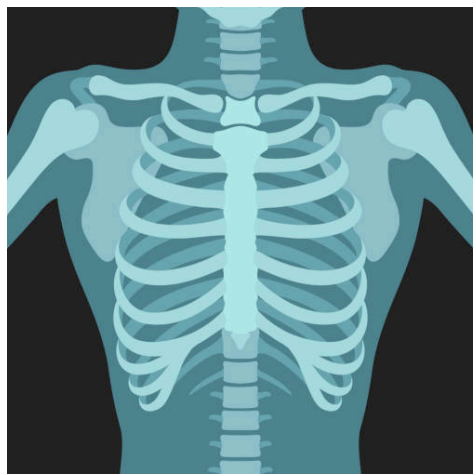


experts

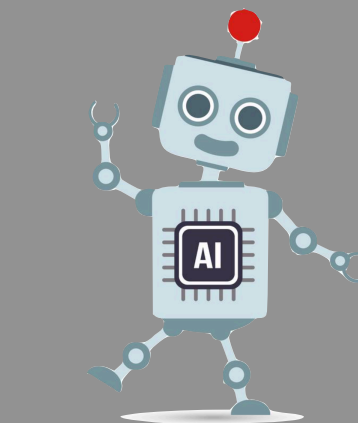
?

expert

input  
features

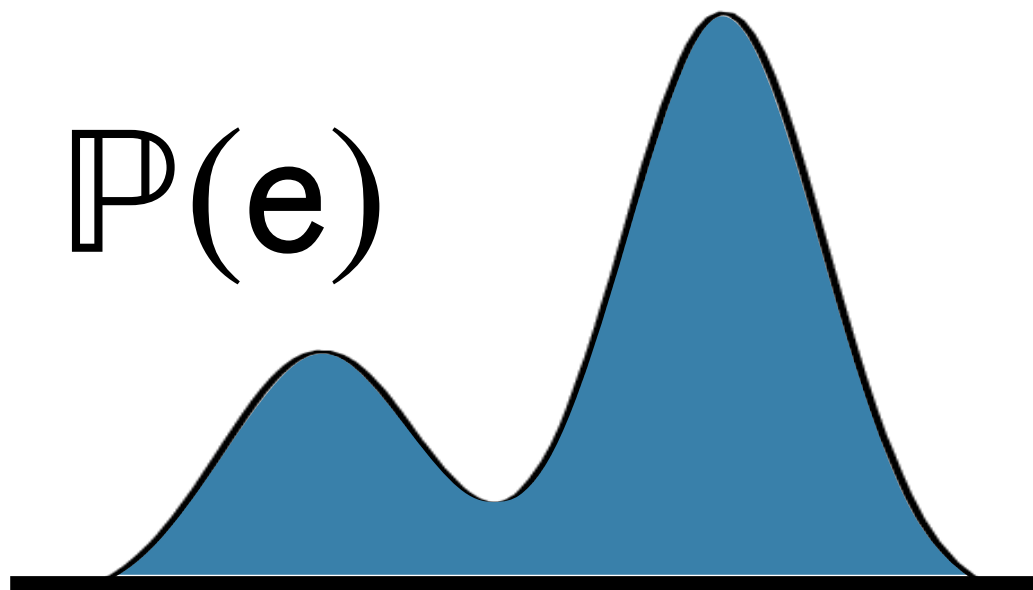


allocation  
mechanism



classifier

$\mathbb{P}(e)$



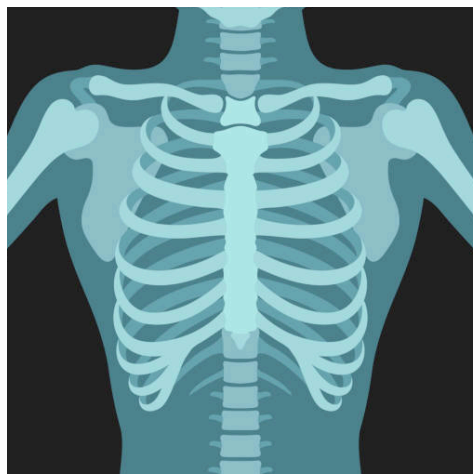
experts

$e \sim \mathbb{P}(e)$

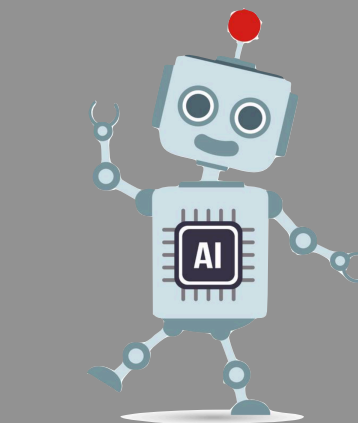


expert

input  
features

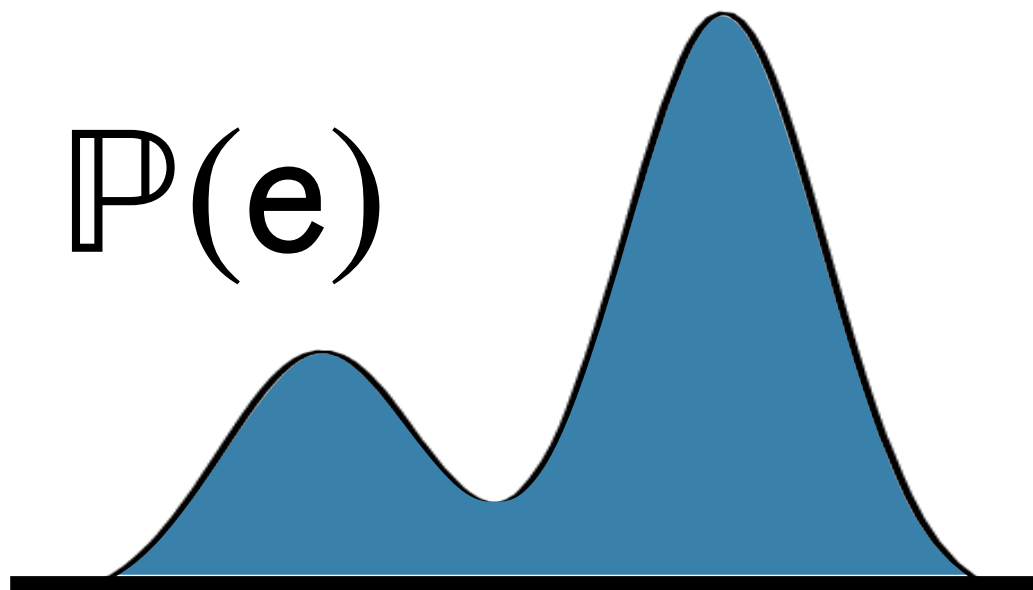


allocation  
mechanism



classifier

$\mathbb{P}(e)$



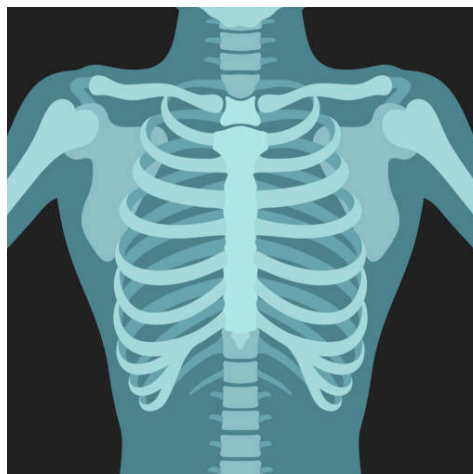
experts

$e \sim \mathbb{P}(e)$

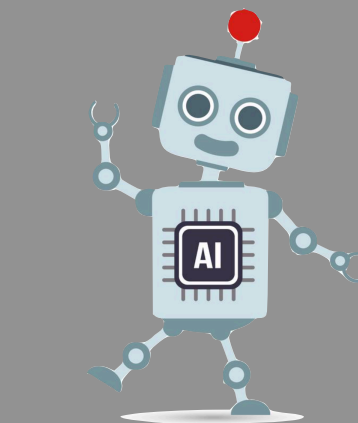


expert

input  
features



allocation  
mechanism



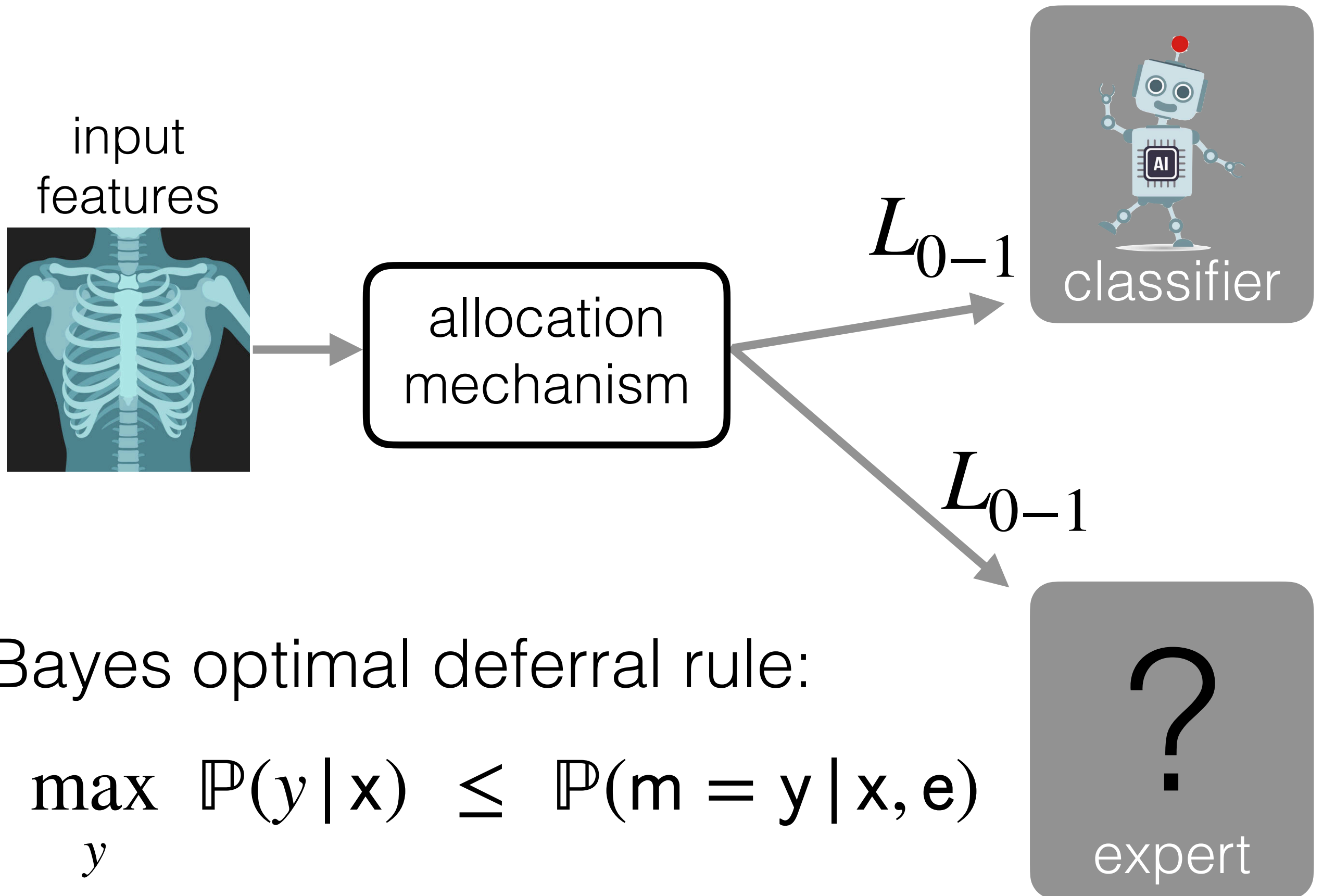
classifier

$$e \sim \mathbb{P}(e)$$
$$m \sim \mathbb{P}(m | e)$$

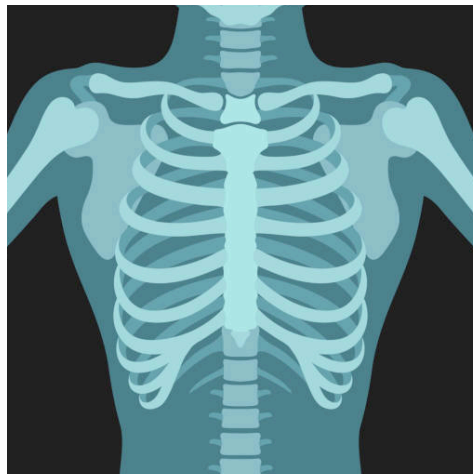


expert

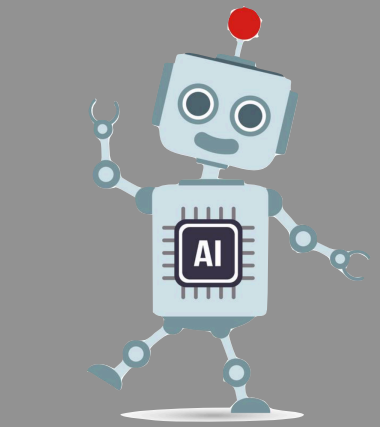




input  
features



allocation  
mechanism



classifier

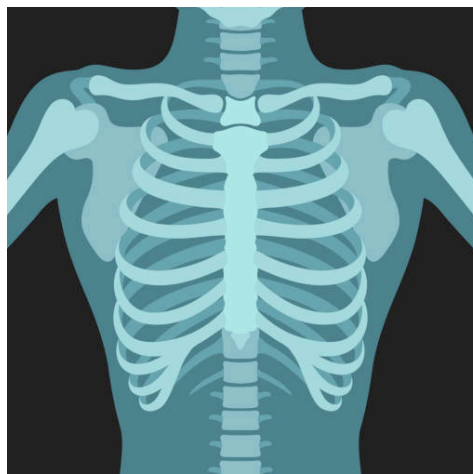
defer to expert if...

$$\max_{y \in [1, K]} h_y(\mathbf{x}) \leq h_{\perp}(\mathbf{x}, \mathbf{e})$$

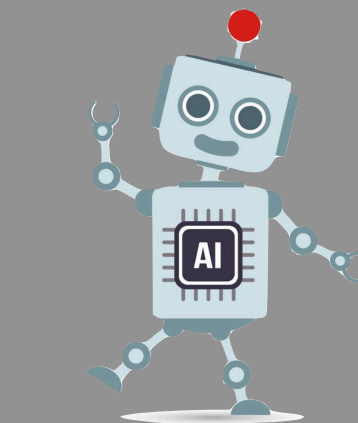
?

expert

input  
features



allocation  
mechanism



classifier

defer to expert if...

$$\max_{y \in [1, K]} h_y(\mathbf{x}) \leq h_{\perp}(\mathbf{x}, \mathbf{e})$$



expert

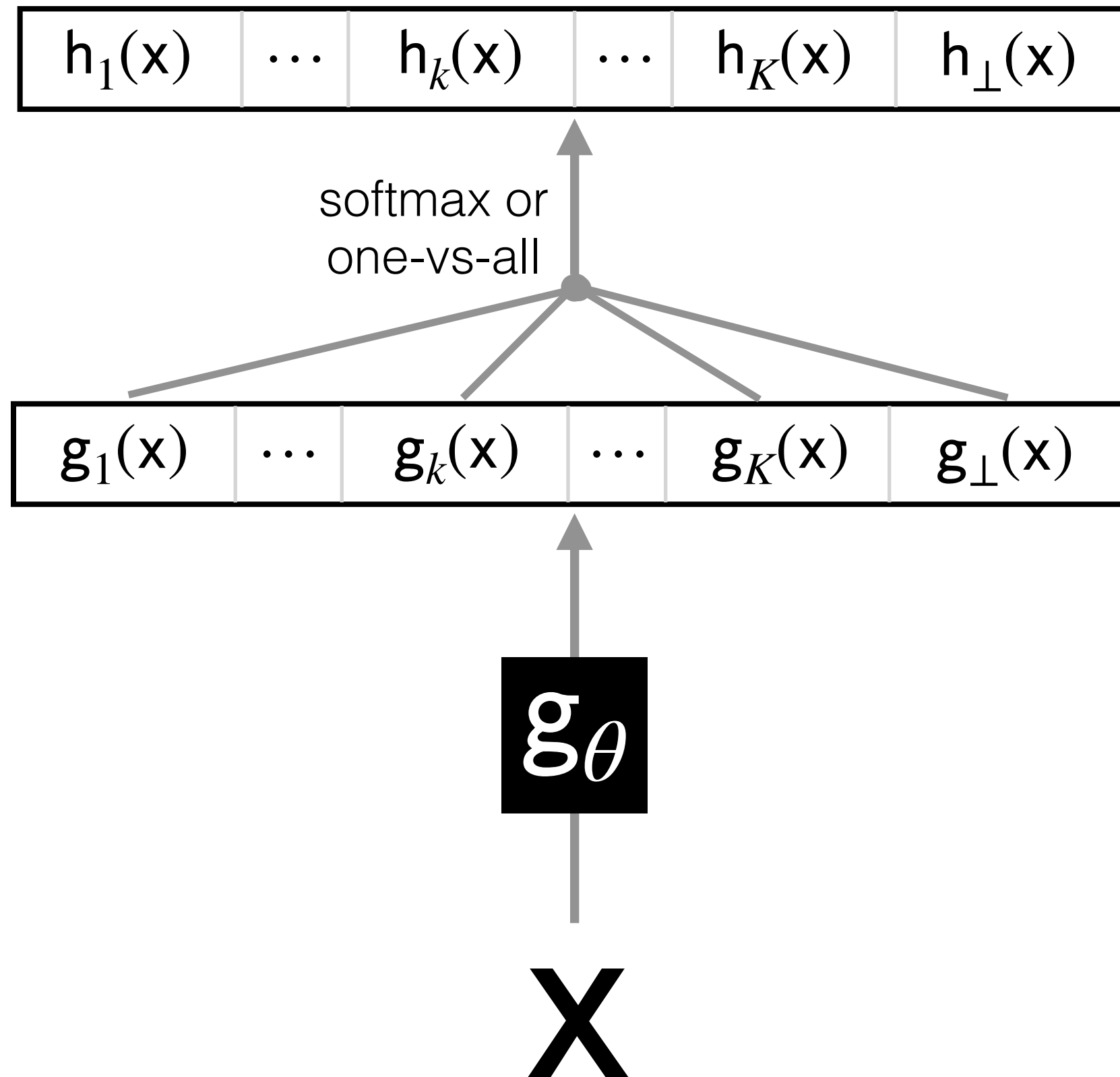
- ⊗ single expert
  - ⊗ softmax surrogate loss
  - ⊗ improving calibration via one-vs-all

- ⊗ multiple experts
  - ⊗ surrogate losses
  - ⊗ conformal sets of experts

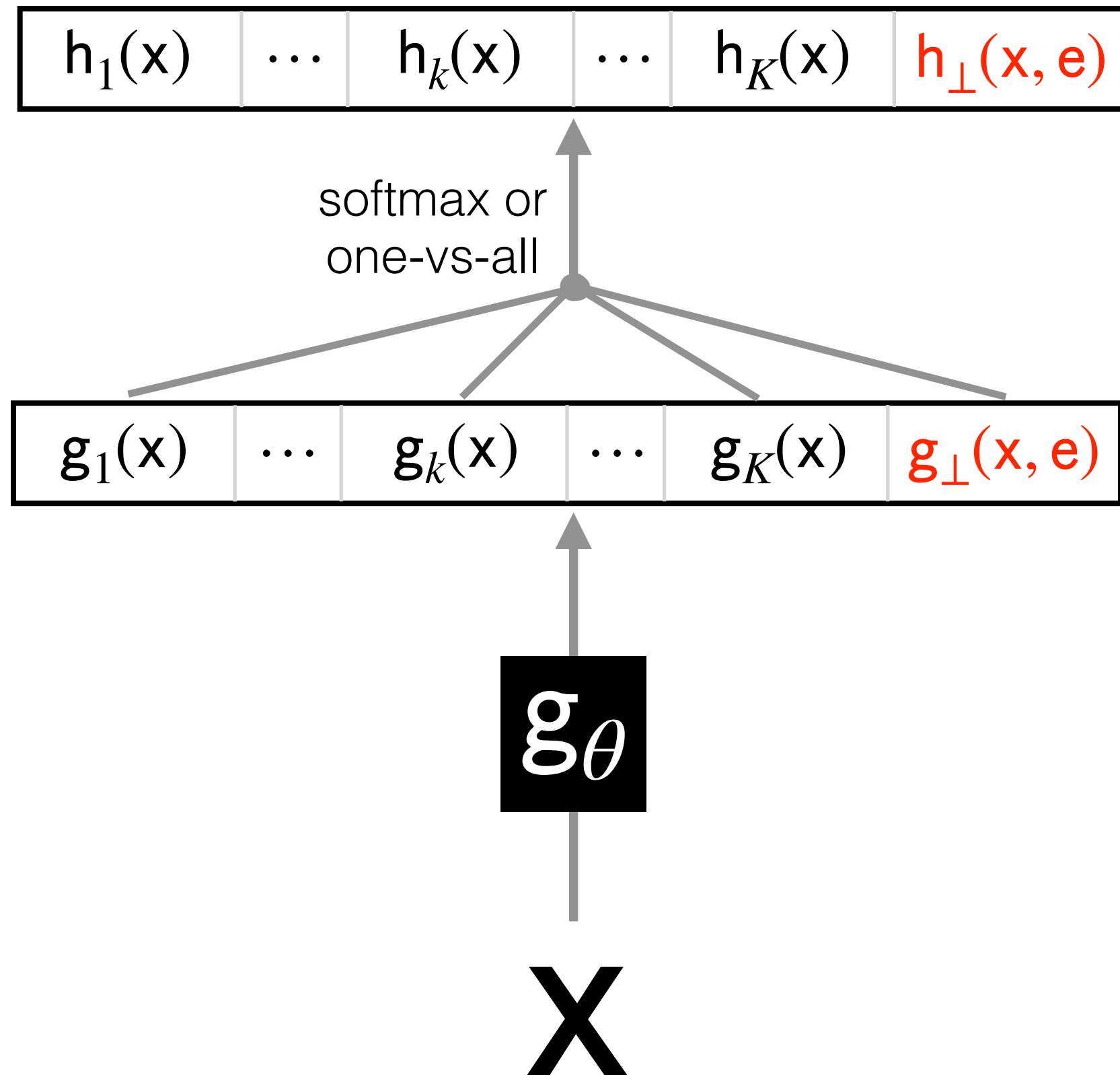
- ⊗ **population of experts**
  - ⊗ surrogate losses
  - ⊗ meta-learning a rejector

- ⊗ single expert
  - ⊗ softmax surrogate loss
  - ⊗ improving calibration via one-vs-all
- ⊗ multiple experts
  - ⊗ surrogate losses
  - ⊗ conformal sets of experts
- ⊗ **population of experts**
  - ⊗ surrogate losses
  - ⊗ meta-learning a rejector

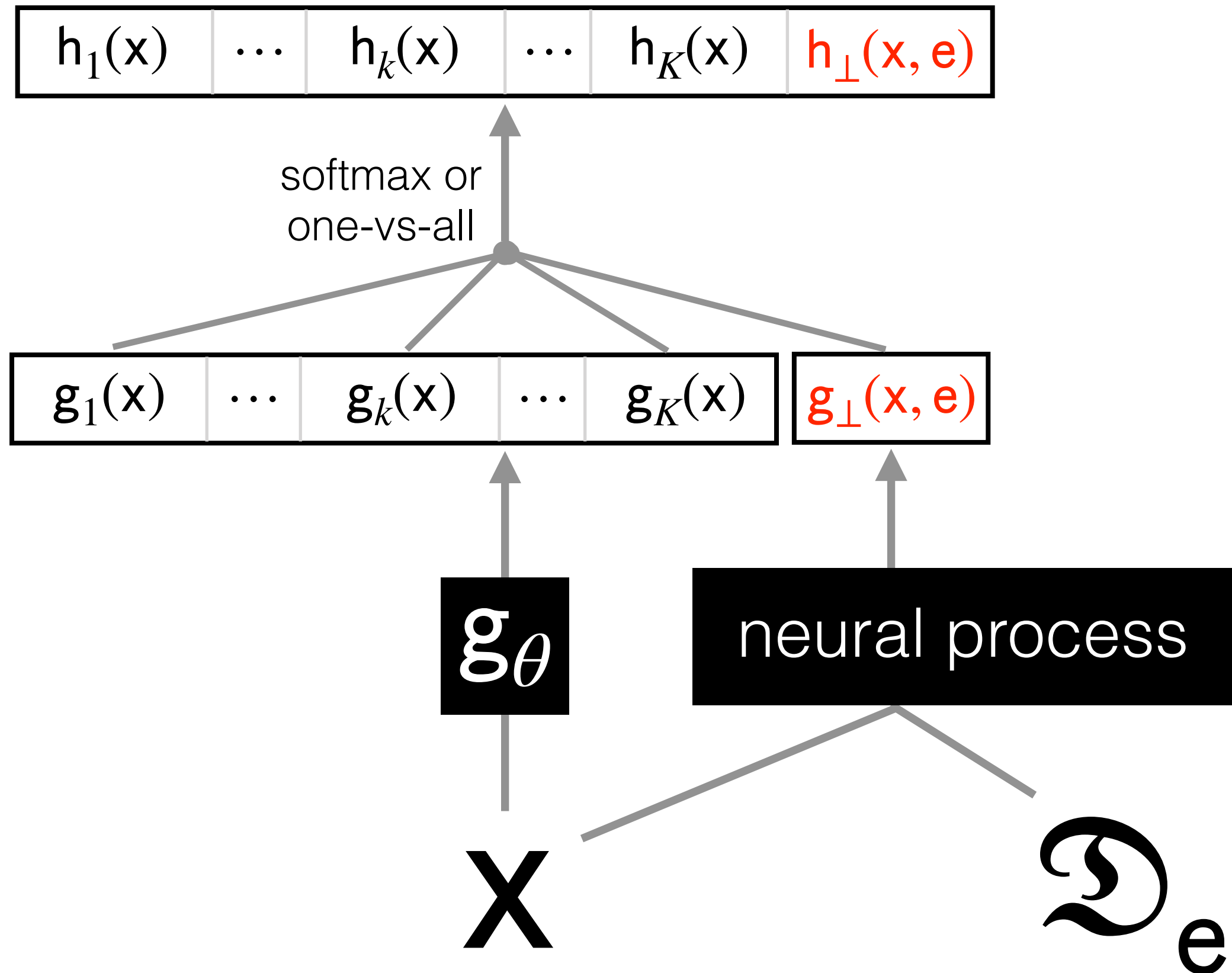
# meta-learning implementation



# meta-learning implementation

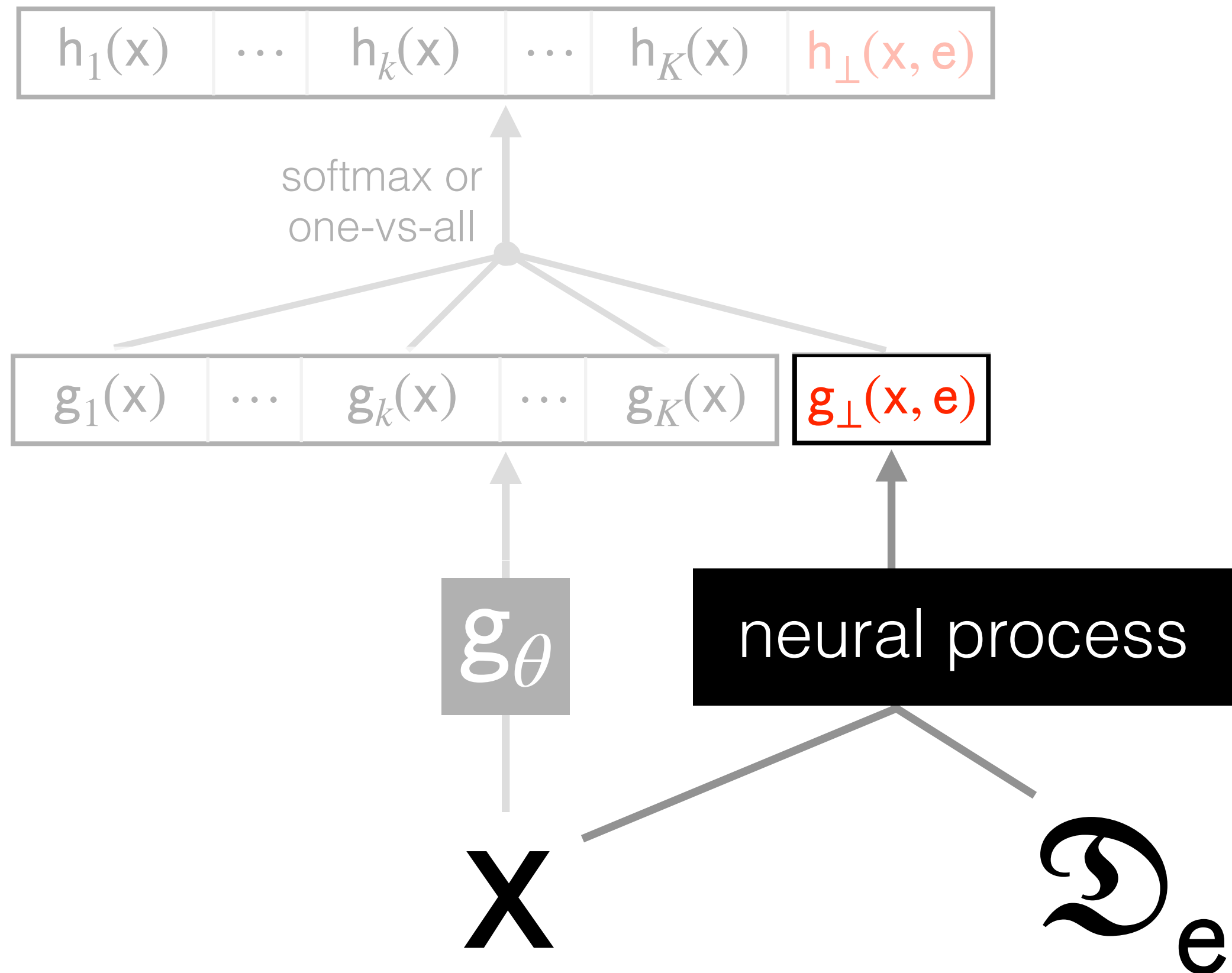


# meta-learning implementation

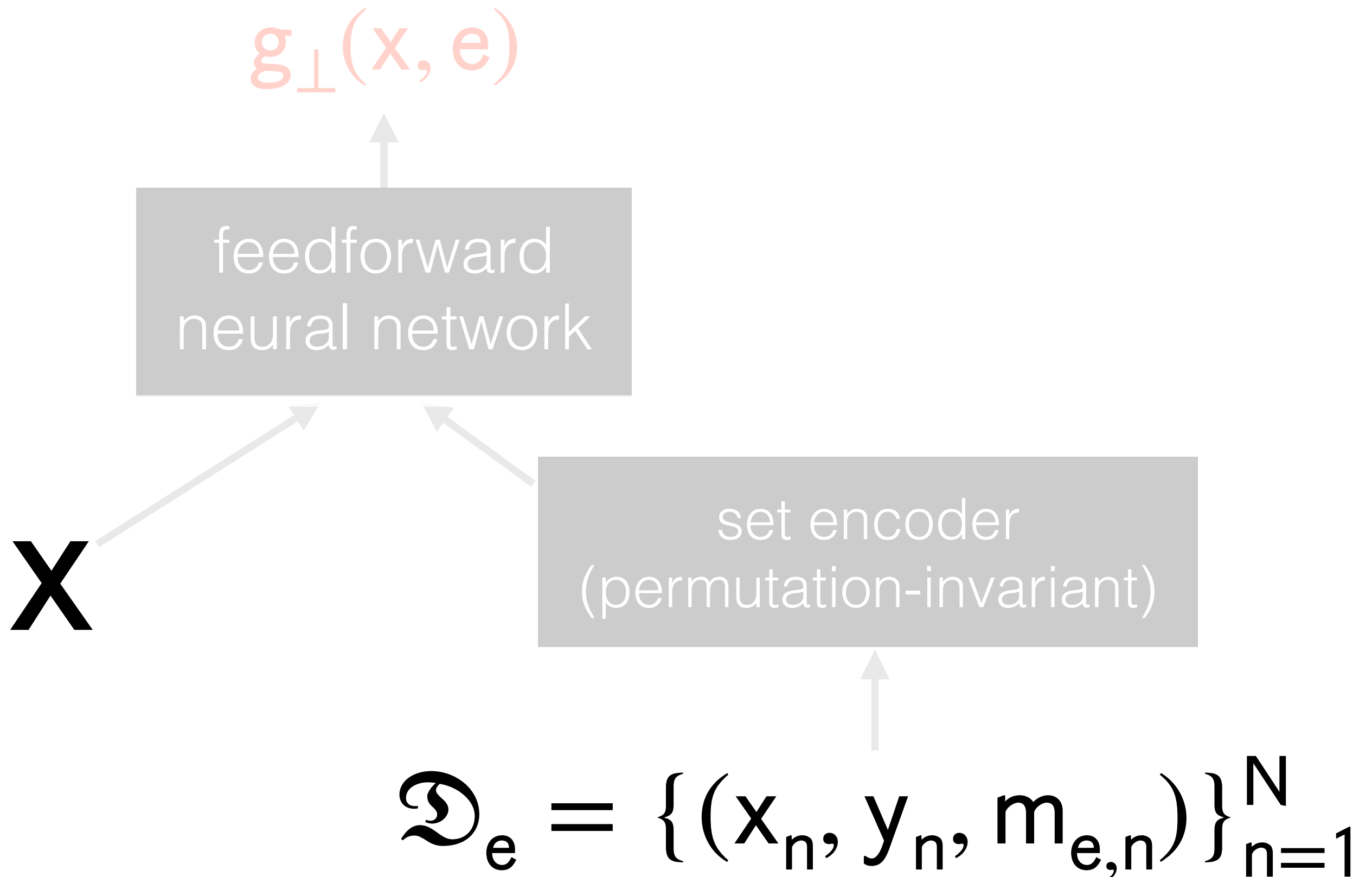




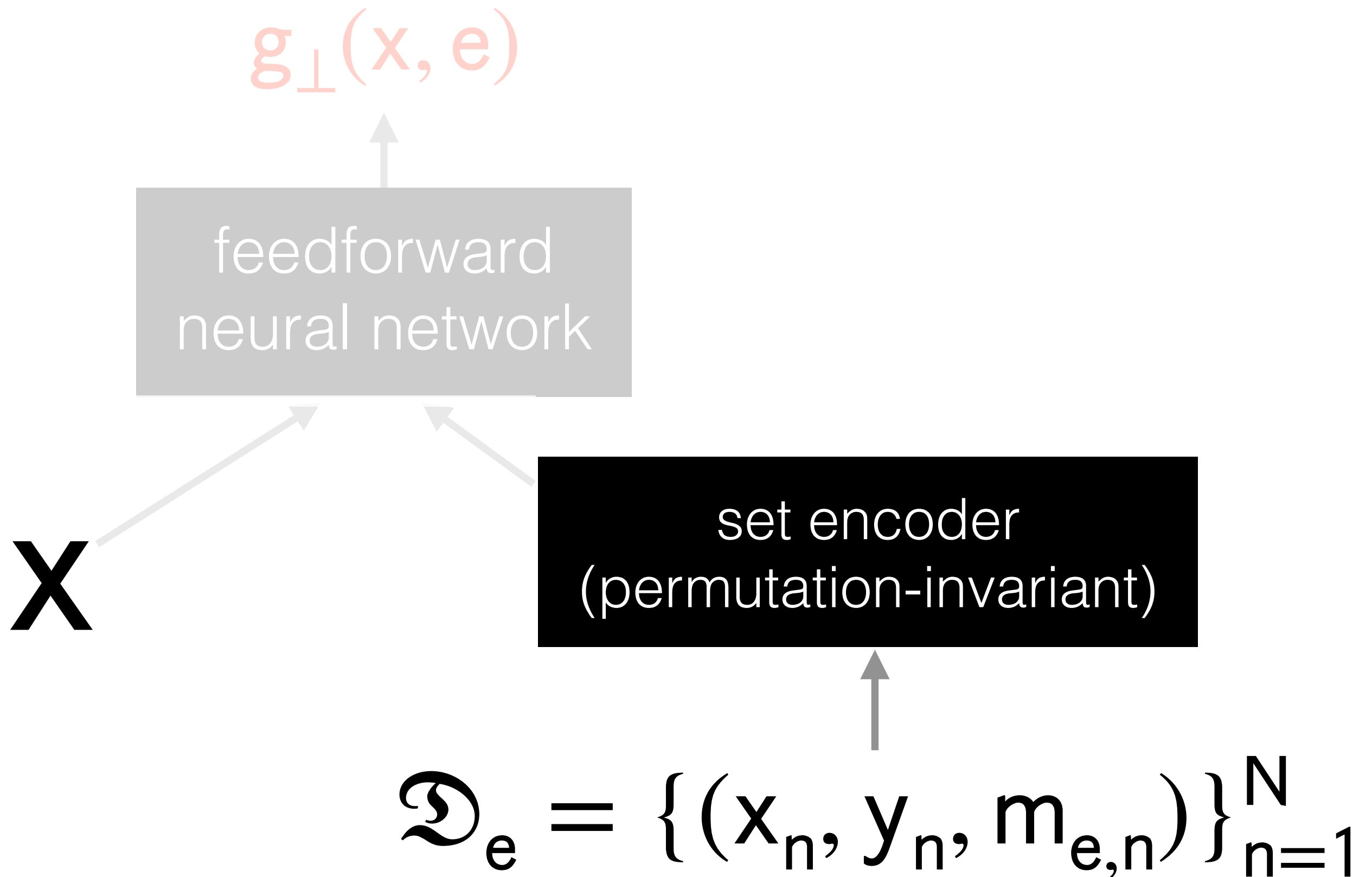
# meta-learning implementation



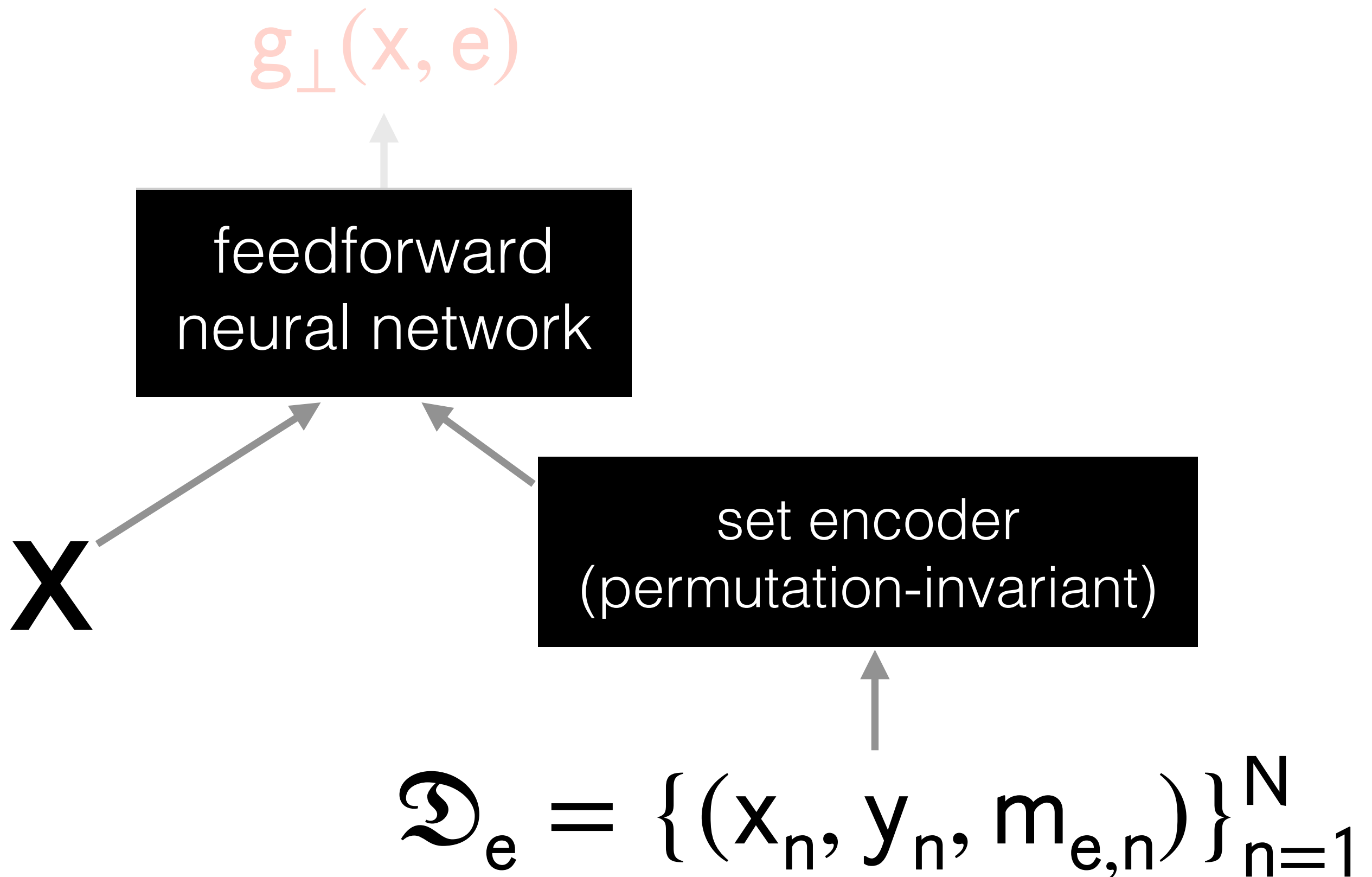
# neural process rejector



# neural process rejector

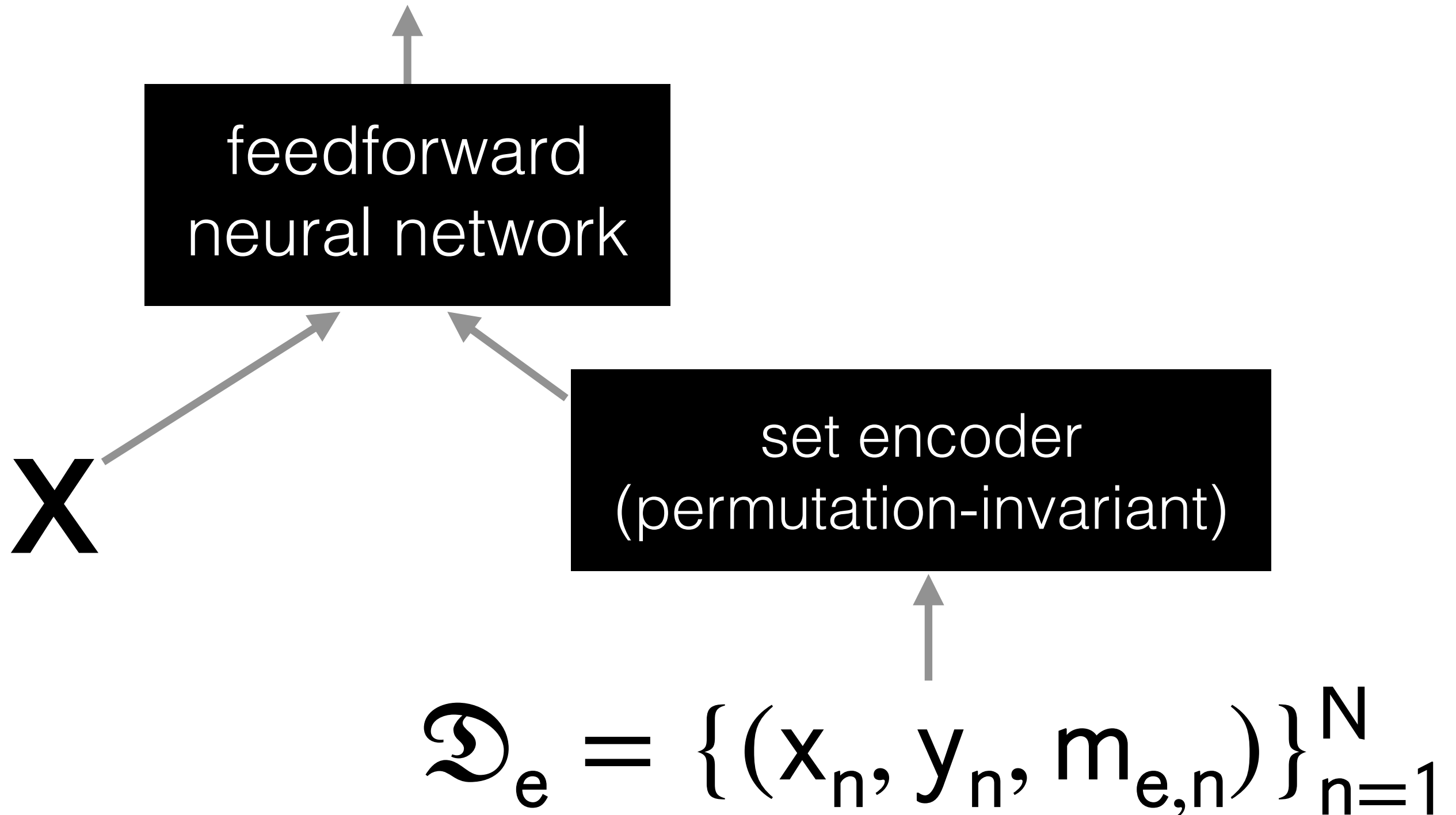


# neural process rejector



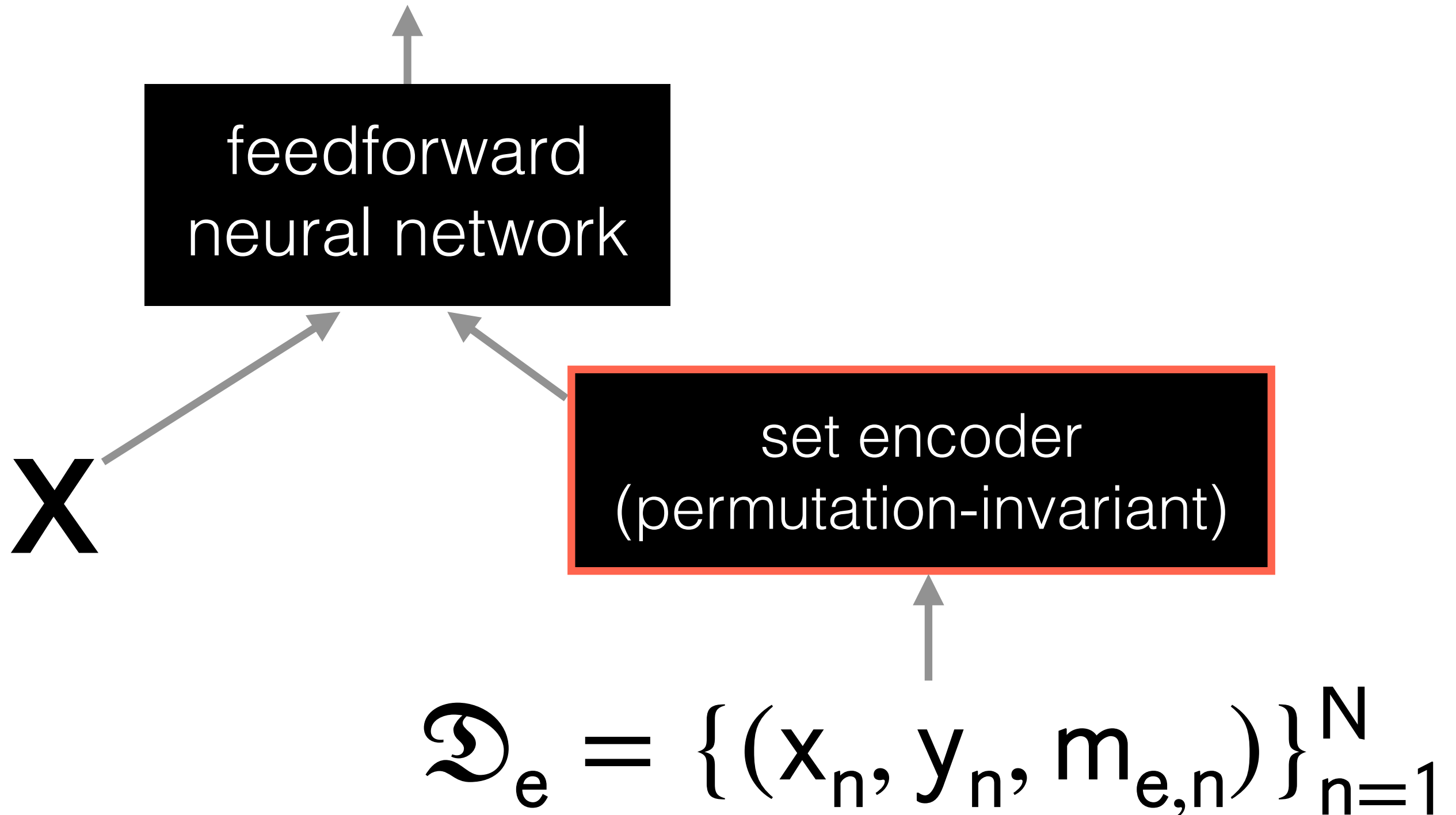
# neural process rejector

$$g_{\perp}(x, e)$$

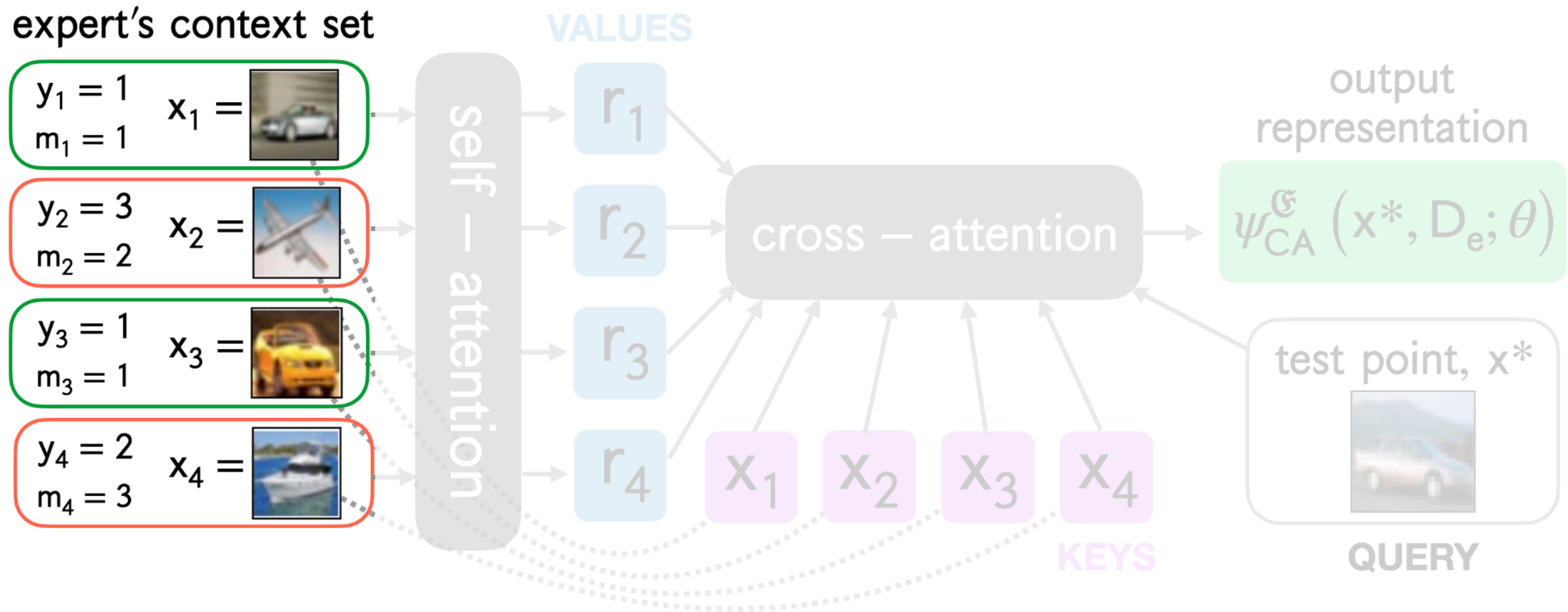


# neural process rejector

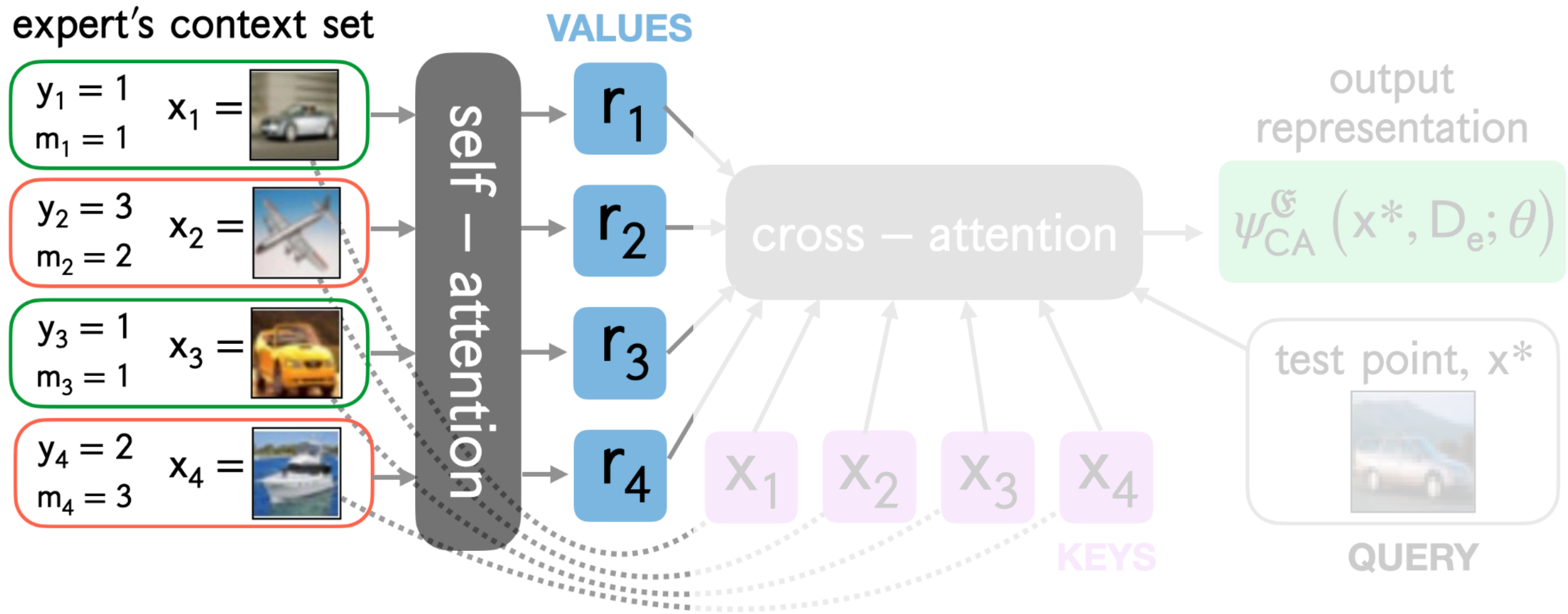
$$g_{\perp}(x, e)$$



# neural process rejector

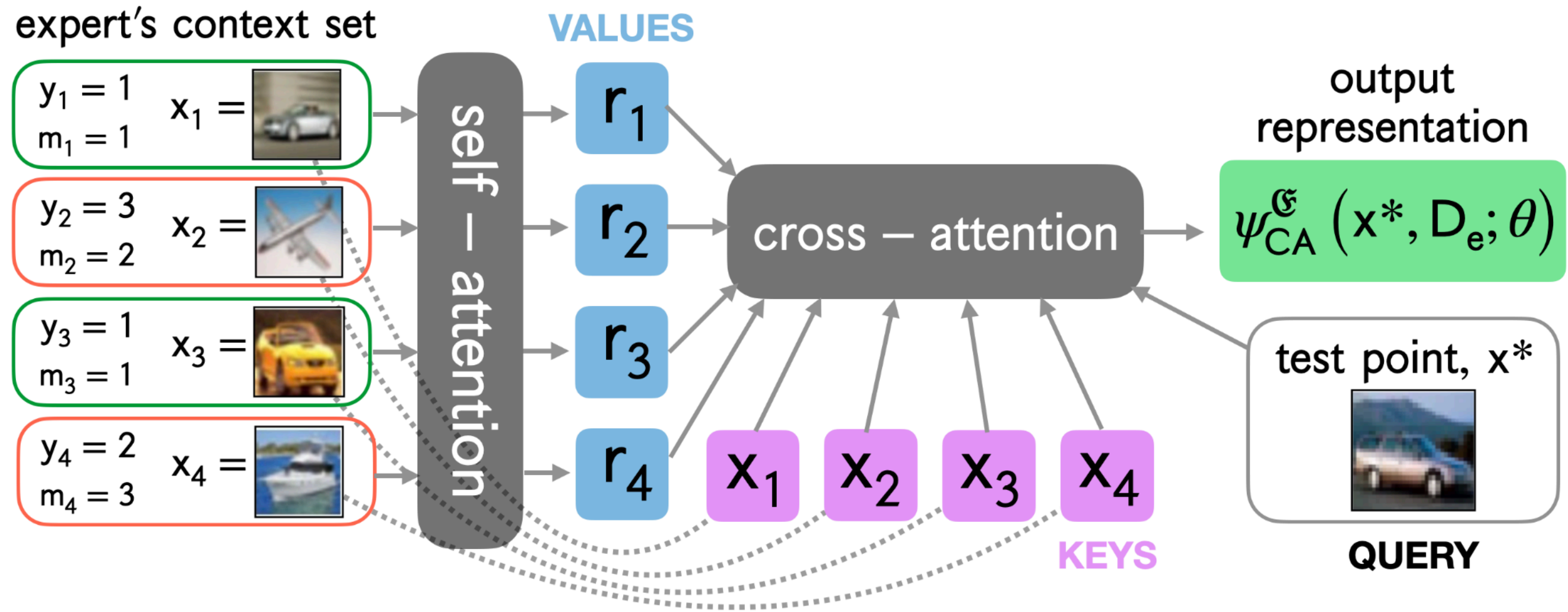


# neural process rejector



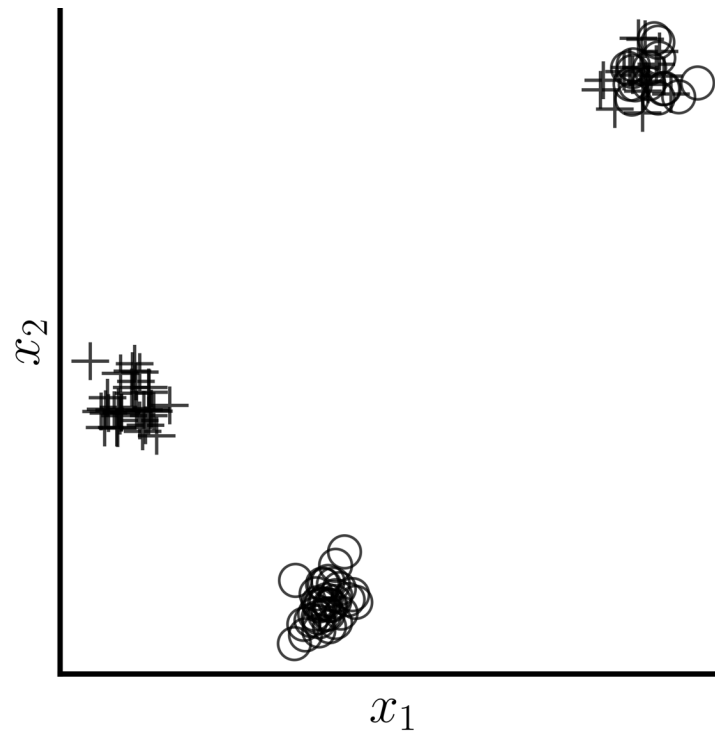


# neural process rejector



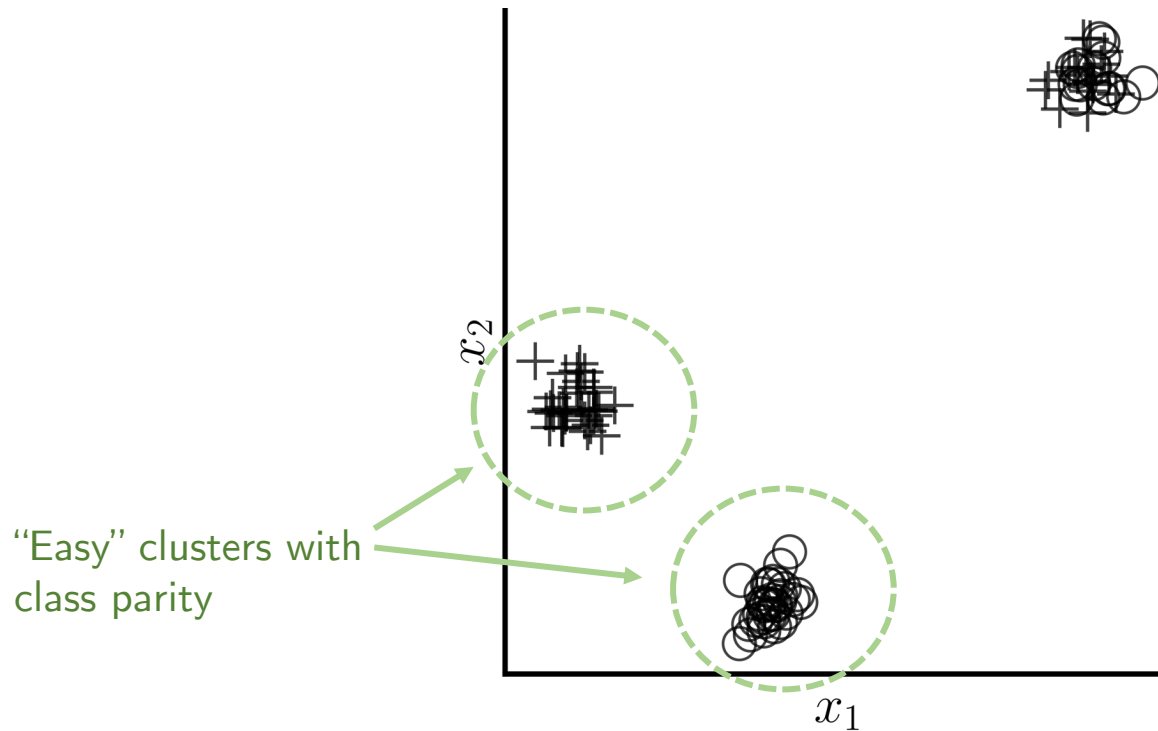
# Experiment: Synthetic data

+ : class 0    O : class 1



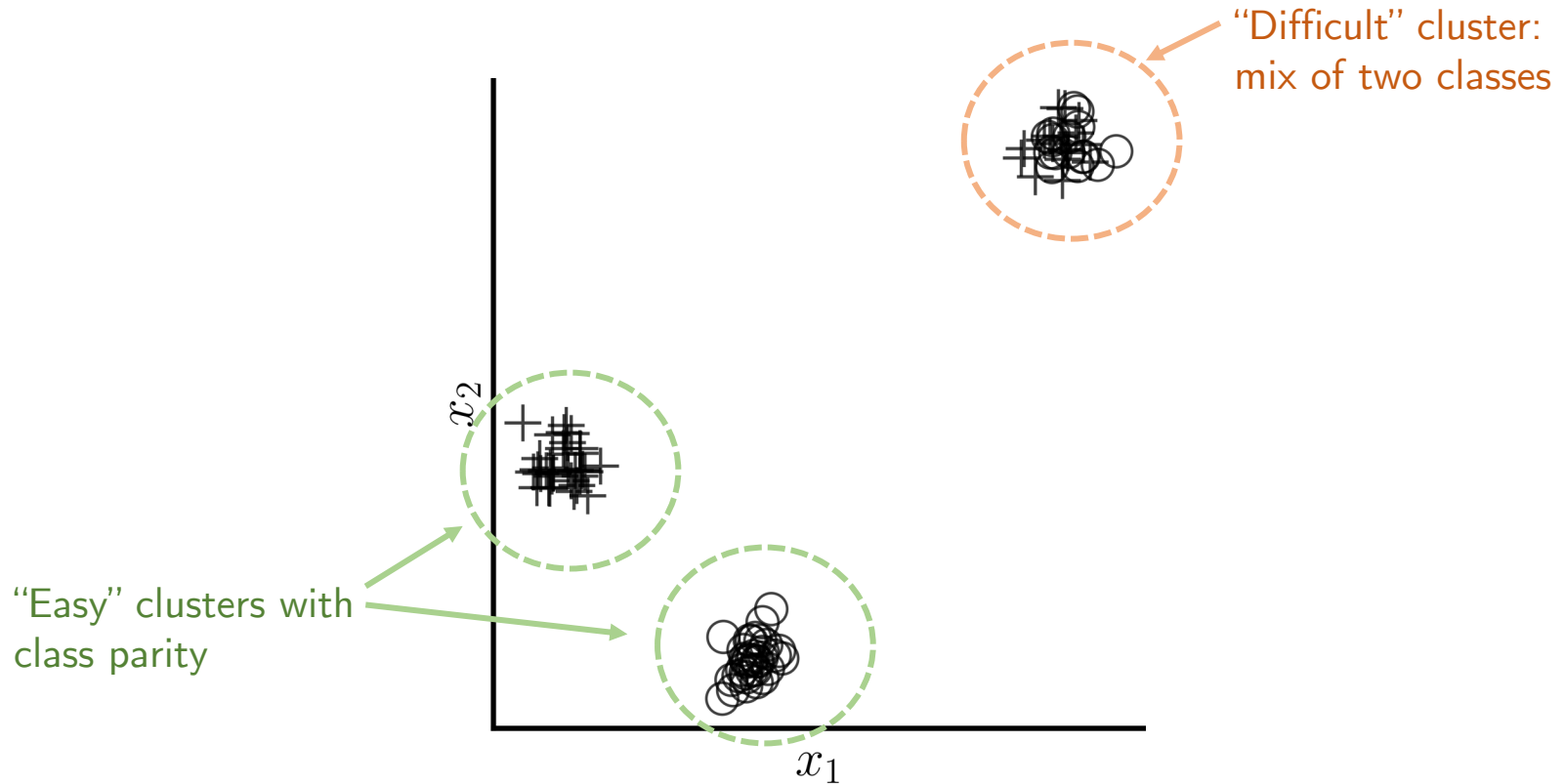
# Experiment: Synthetic data

+ : class 0    O : class 1





# Experiment: Synthetic data

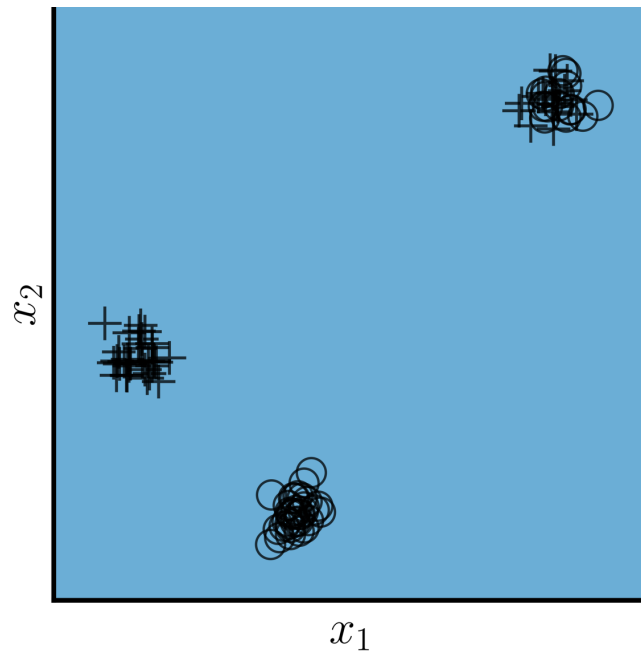
+ : class 0    O : class 1



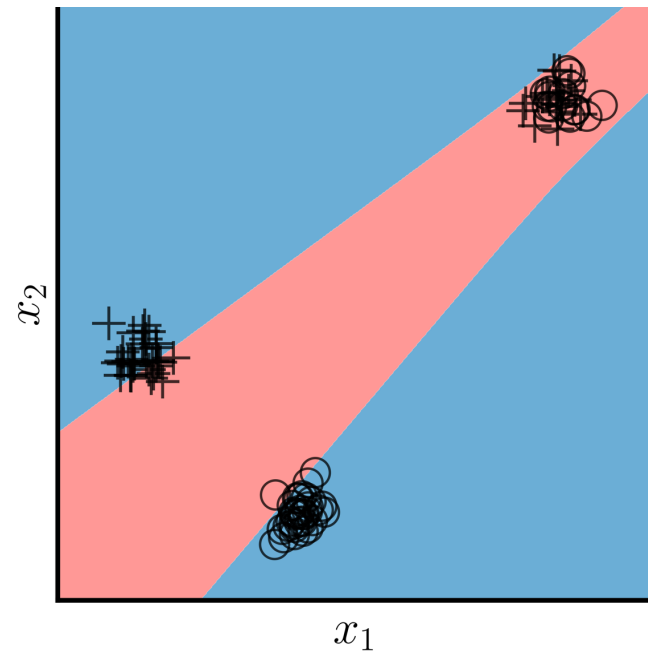
# Experiment: Synthetic data

$+$  : class 0     $\circ$  : class 1     L2D-Pop classifier region     L2D-Pop deferral region



Unskilled expert (1% accuracy)



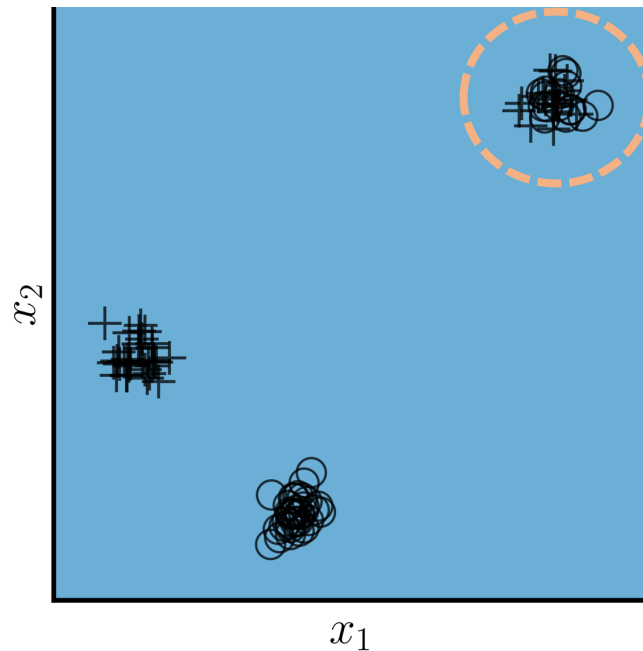
Skilled expert (95% accuracy)



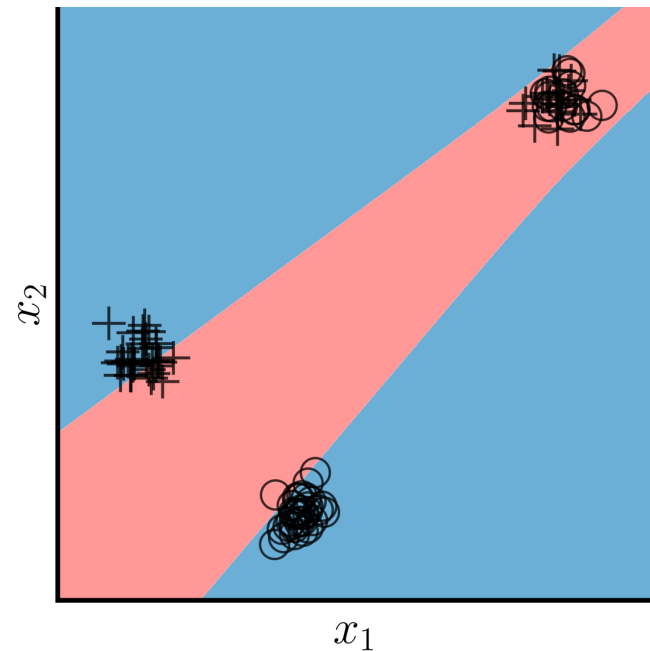
# Experiment: Synthetic data


$+$  : class 0     $O$  : class 1     L2D-Pop classifier region     L2D-Pop deferral region

Unskilled expert (1% accuracy)





Skilled expert (95% accuracy)

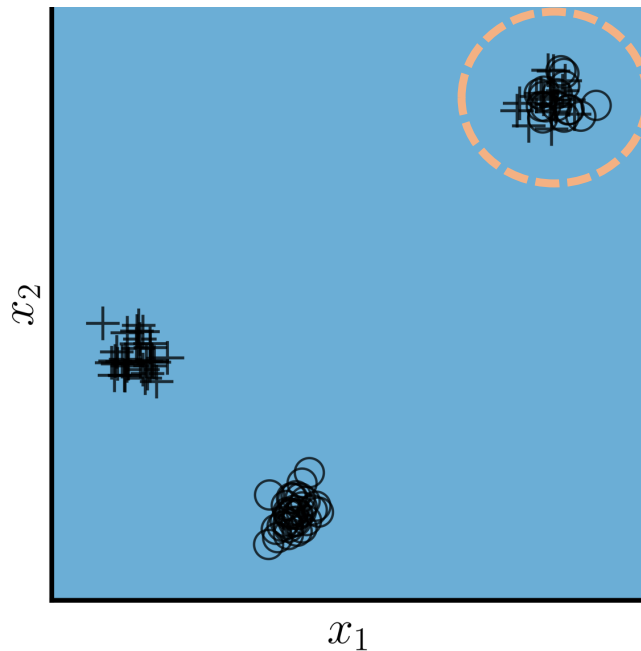



L2D-Pop (adaptive)     Doesn't defer when the expert is poor

# Experiment: Synthetic data

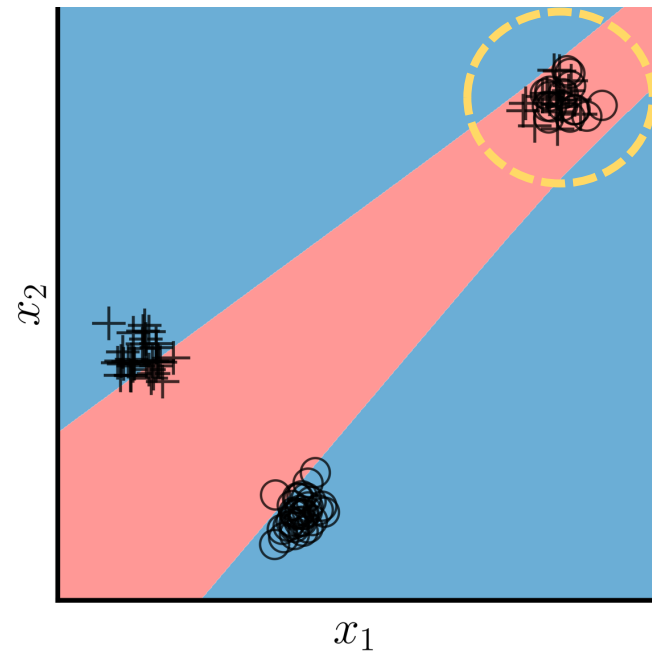
$+$  : class 0     $O$  : class 1     L2D-Pop classifier region     L2D-Pop deferral region


Unskilled expert (1% accuracy)



L2D-Pop (adaptive)     Doesn't defer when the expert is poor

Skilled expert (95% accuracy)

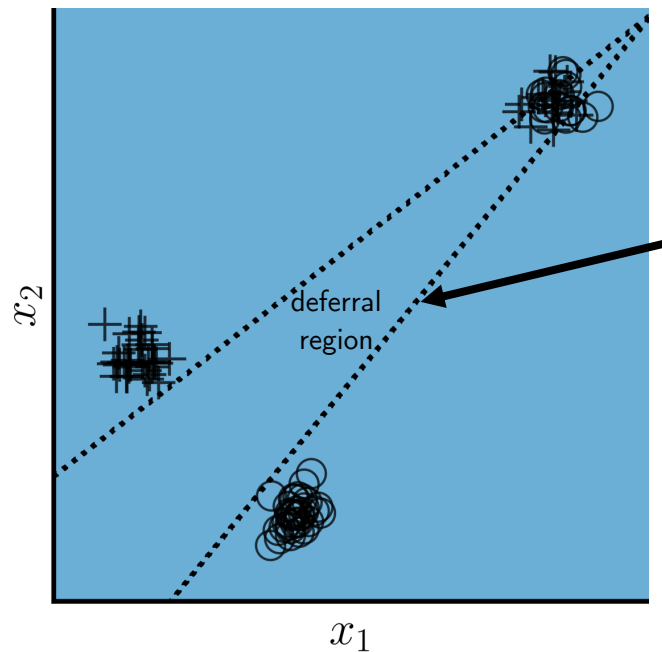


 Defers whole of difficult cluster when expert is good

# Experiment: Synthetic data

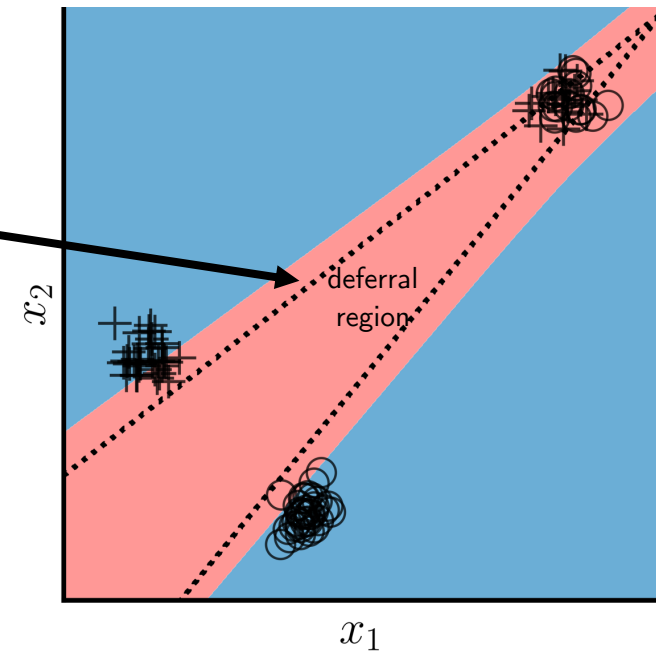
$+$  : class 0     $\circ$  : class 1     L2D-Pop classifier region     L2D-Pop deferral region

Unskilled expert (1% accuracy)



L2D-Pop (adaptive)    ✓ Doesn't defer when the expert is poor

Skilled expert (95% accuracy)



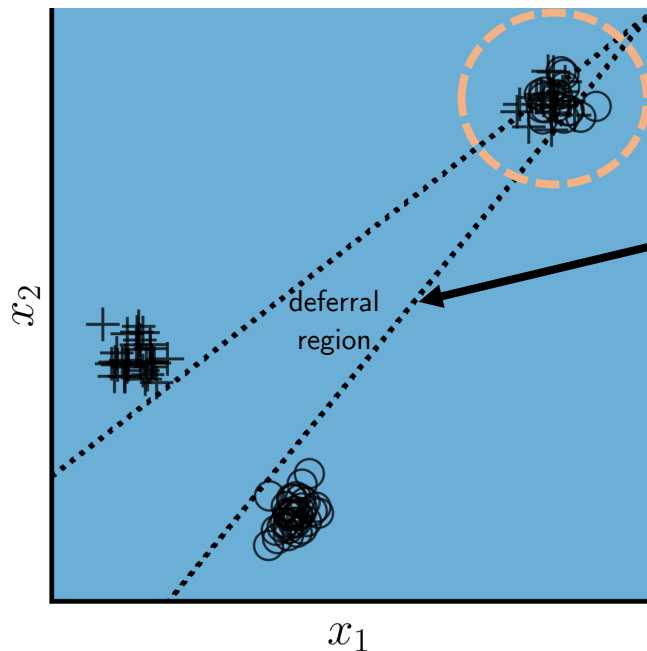
✓ Defers whole of difficult cluster when expert is good



# Experiment: Synthetic data

$+$  : class 0     $O$  : class 1     L2D-Pop classifier region     L2D-Pop deferral region

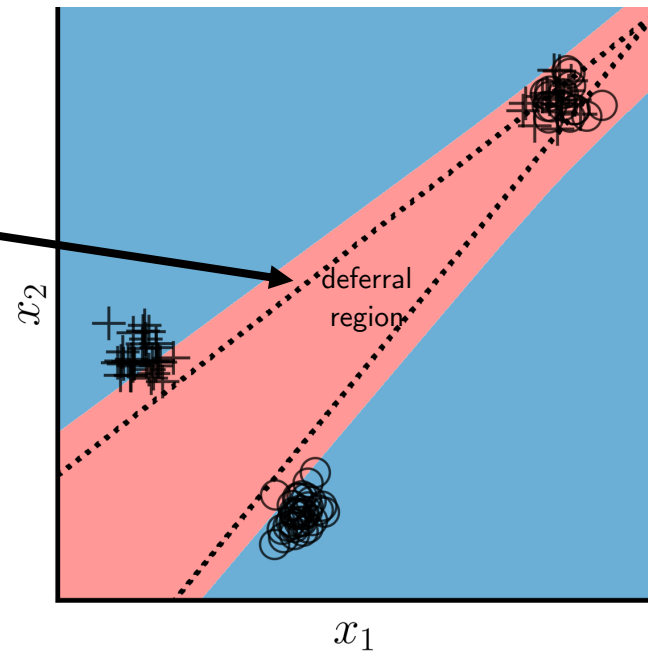
Unskilled expert (1% accuracy)



L2D-Pop (adaptive)    ✓ Doesn't defer when the expert is poor

single-L2D (constant)    ✗ Over-defers as expert does worse than random on difficult cluster

Skilled expert (95% accuracy)

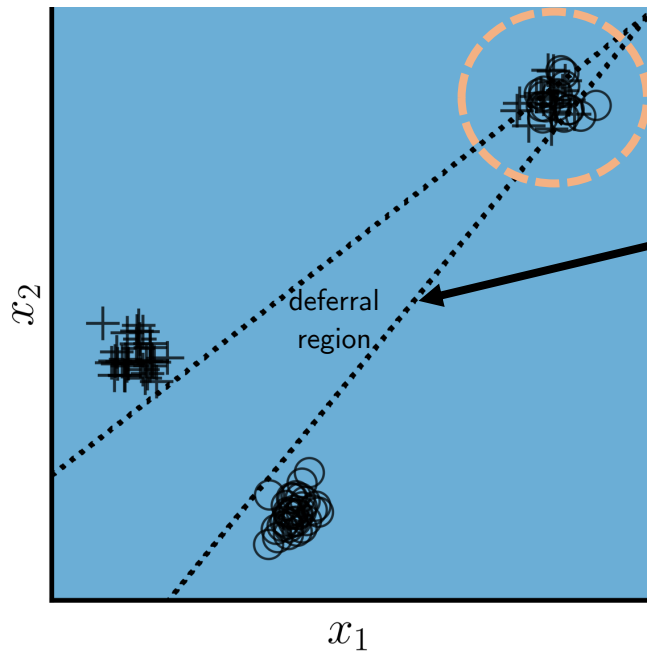


✓ Defers whole of difficult cluster when expert is good

# Experiment: Synthetic data

+ : class 0    O : class 1     L2D-Pop classifier region     L2D-Pop deferral region

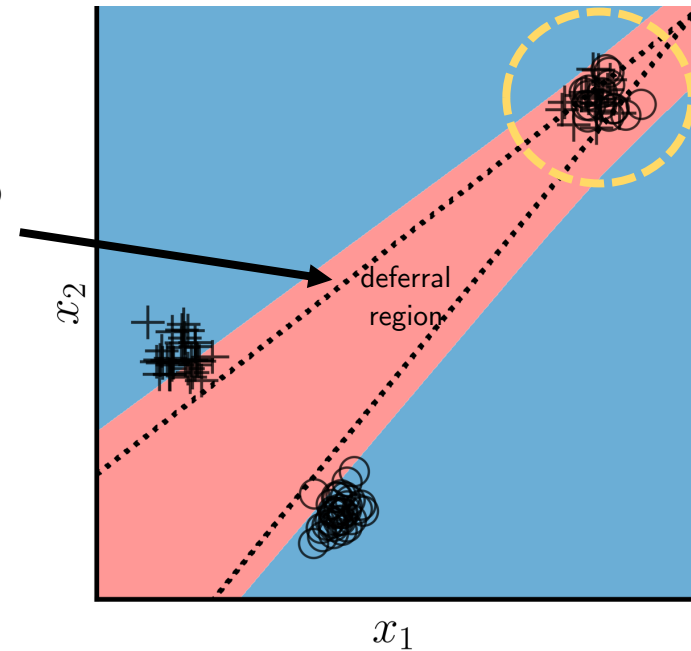
Unskilled expert (1% accuracy)



L2D-Pop (adaptive) ✓ Doesn't defer when the expert is poor

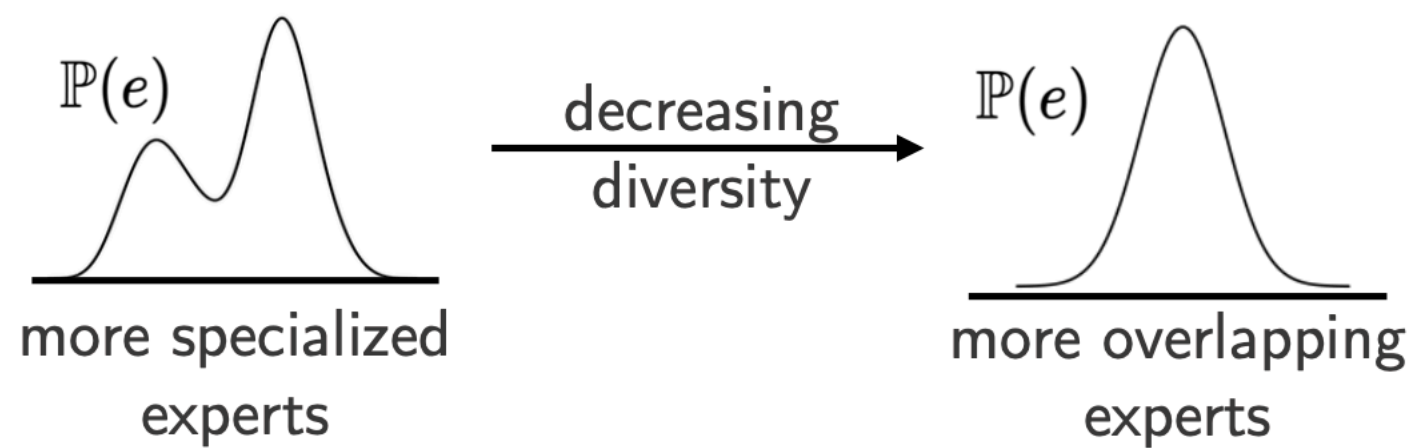
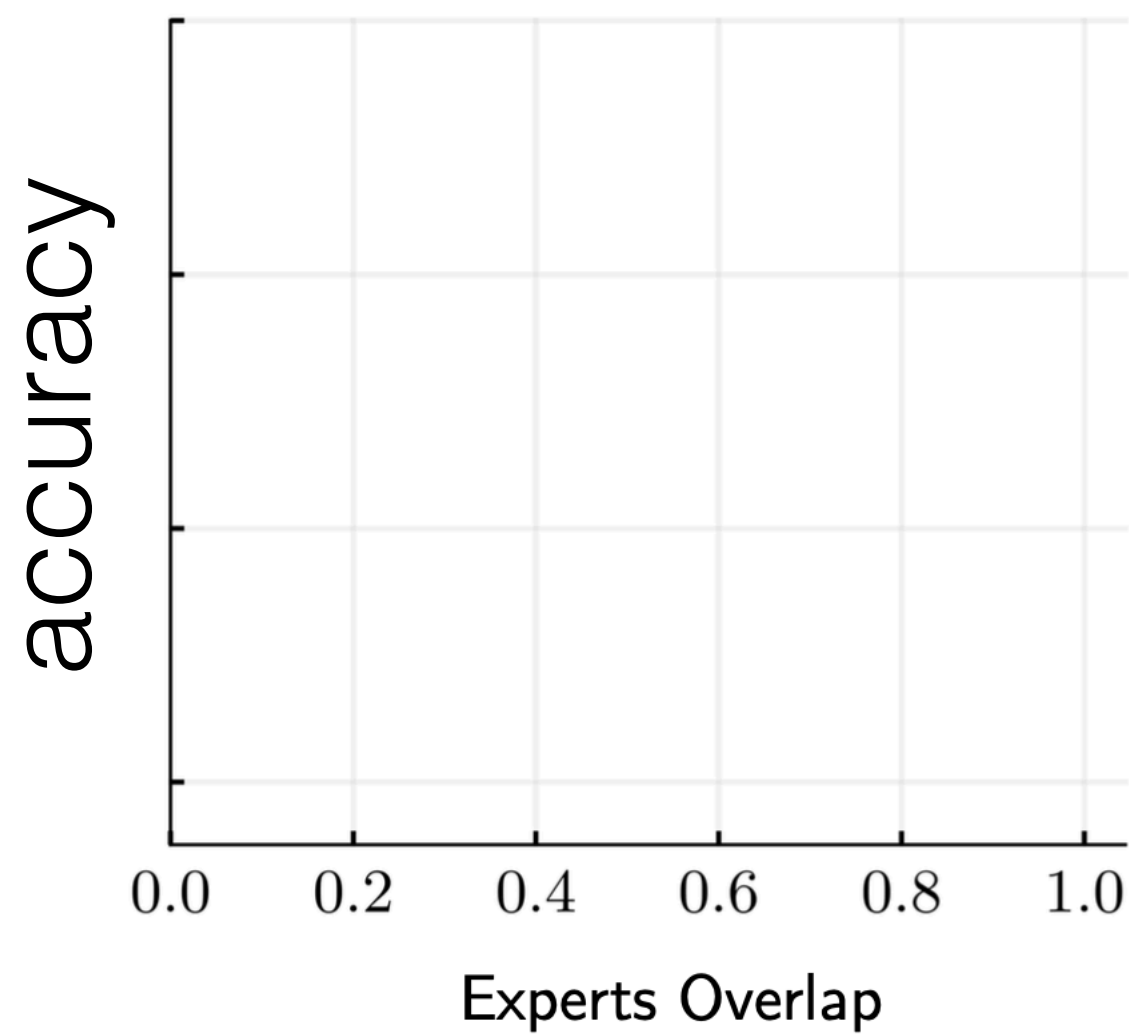
single-L2D (constant) ✗ Over-defers as expert does worse than random on difficult cluster

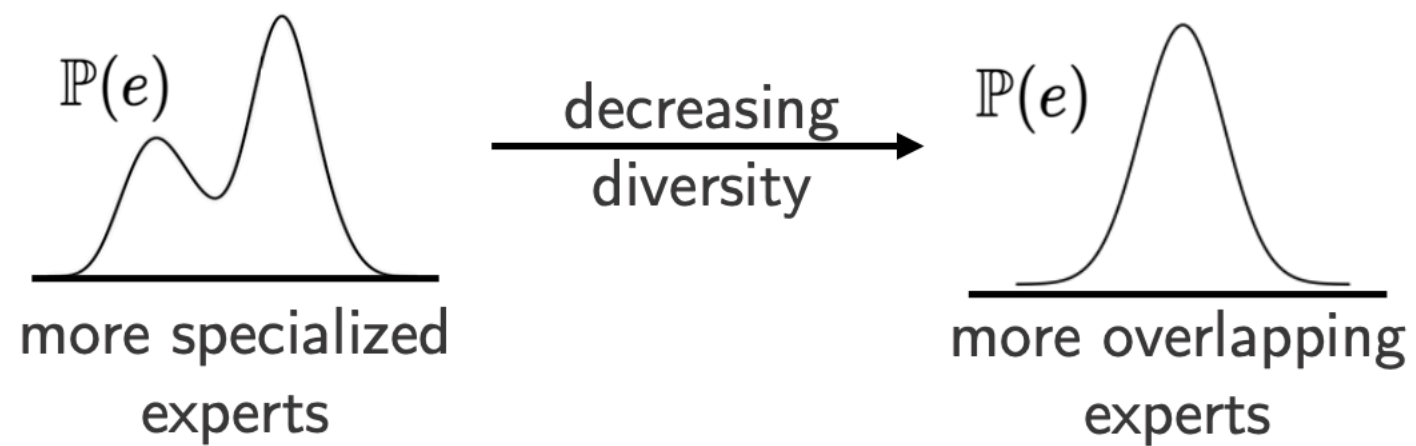
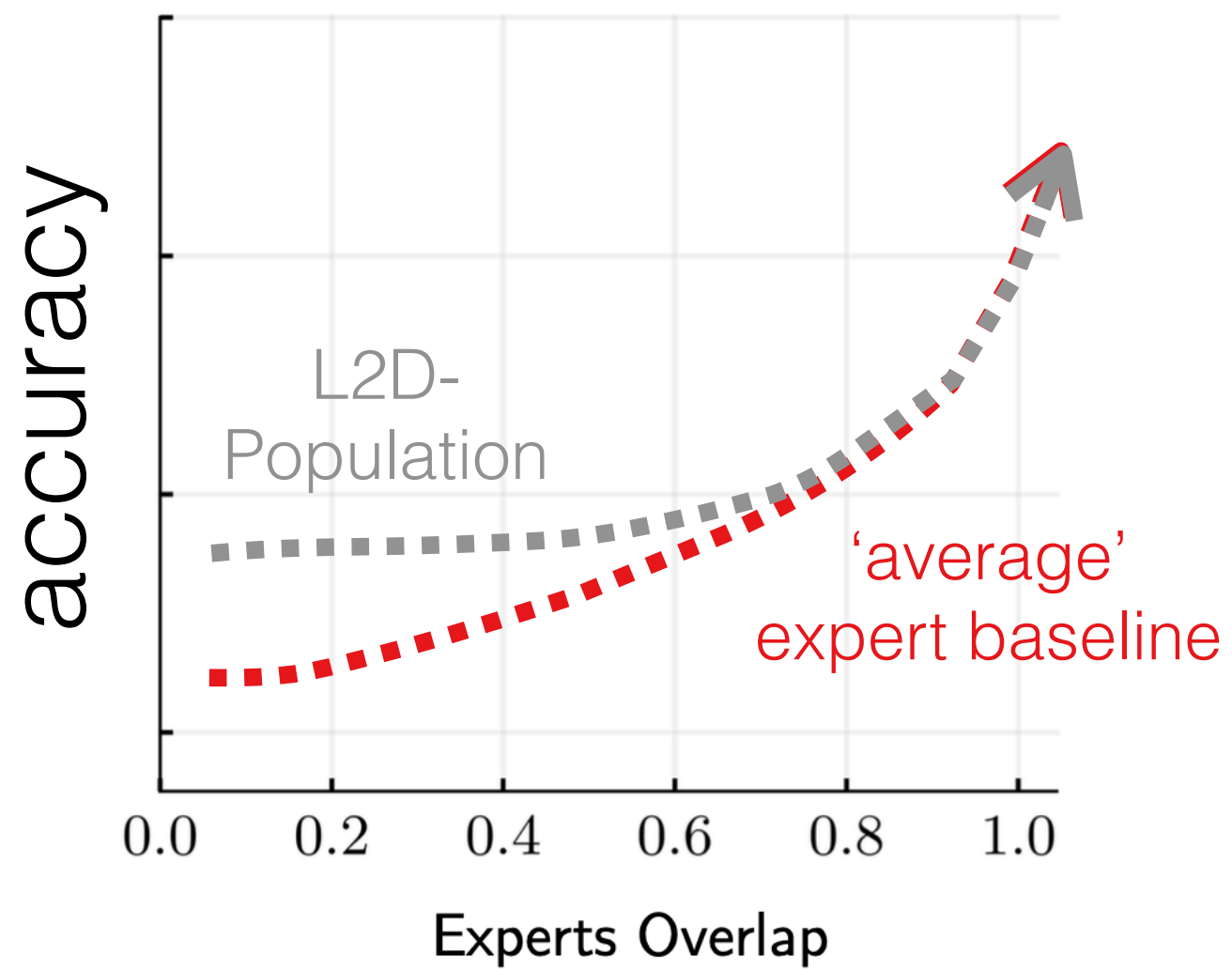
Skilled expert (95% accuracy)

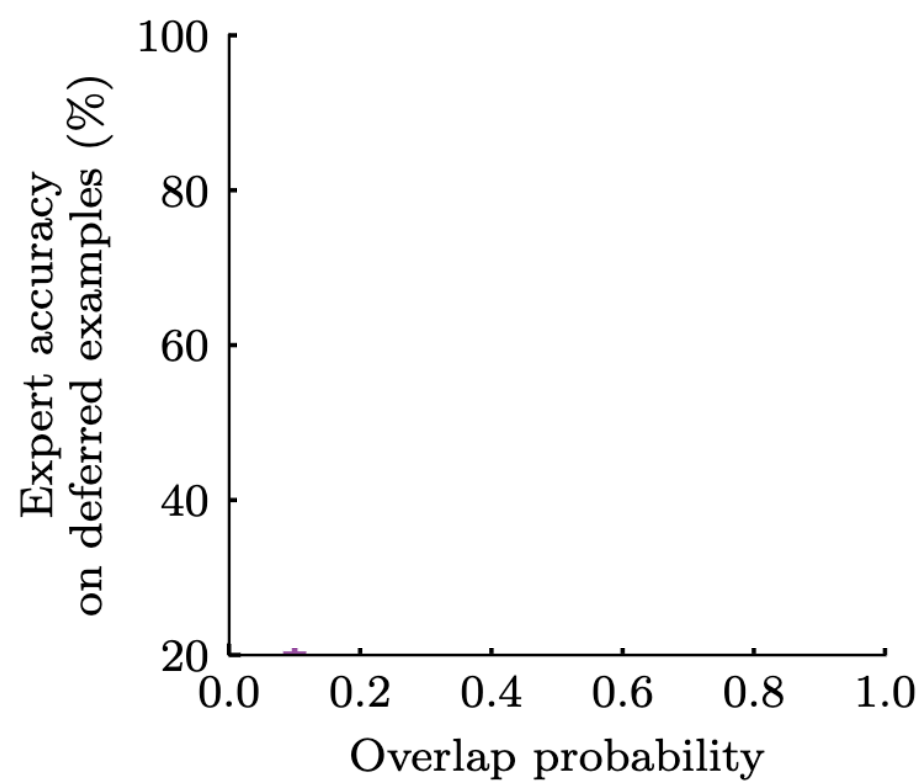
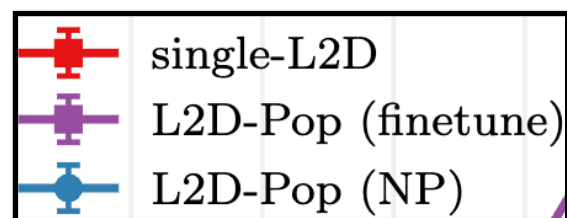


✓ Defers whole of difficult cluster when expert is good

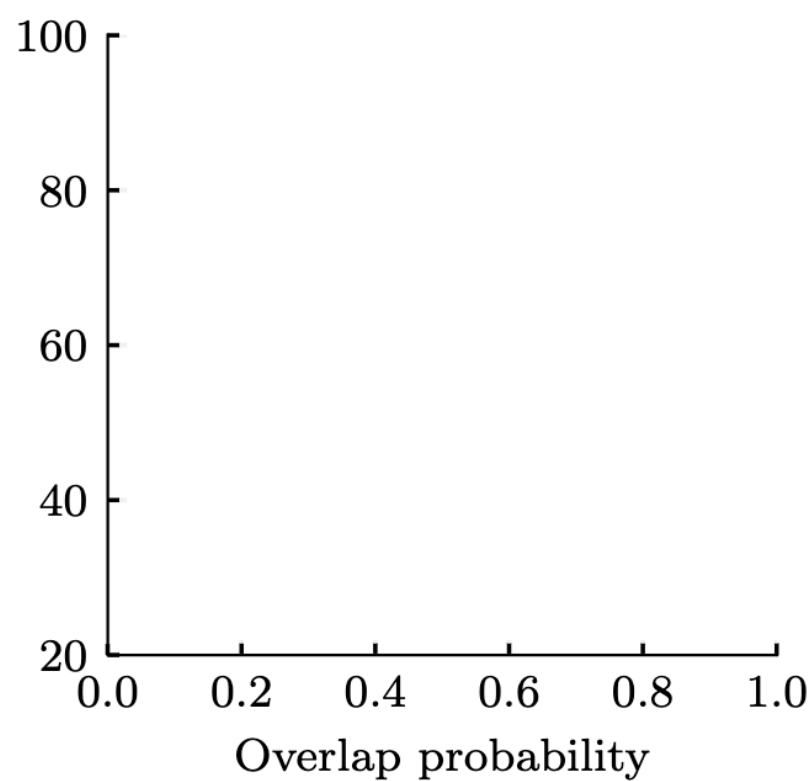
✗ Under-defers as classifier only has random chance of being correct on difficult cluster



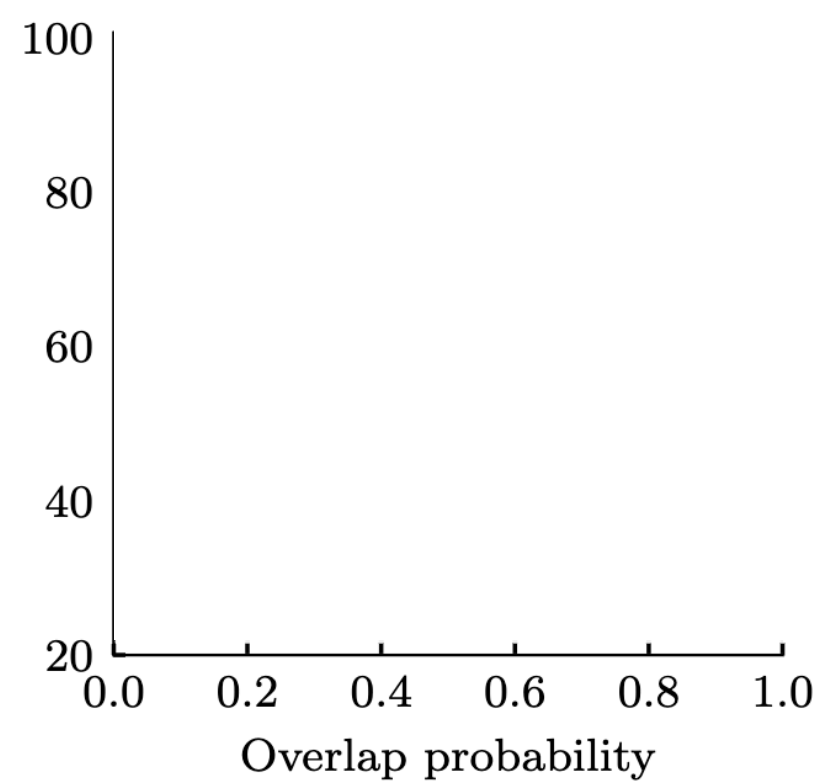




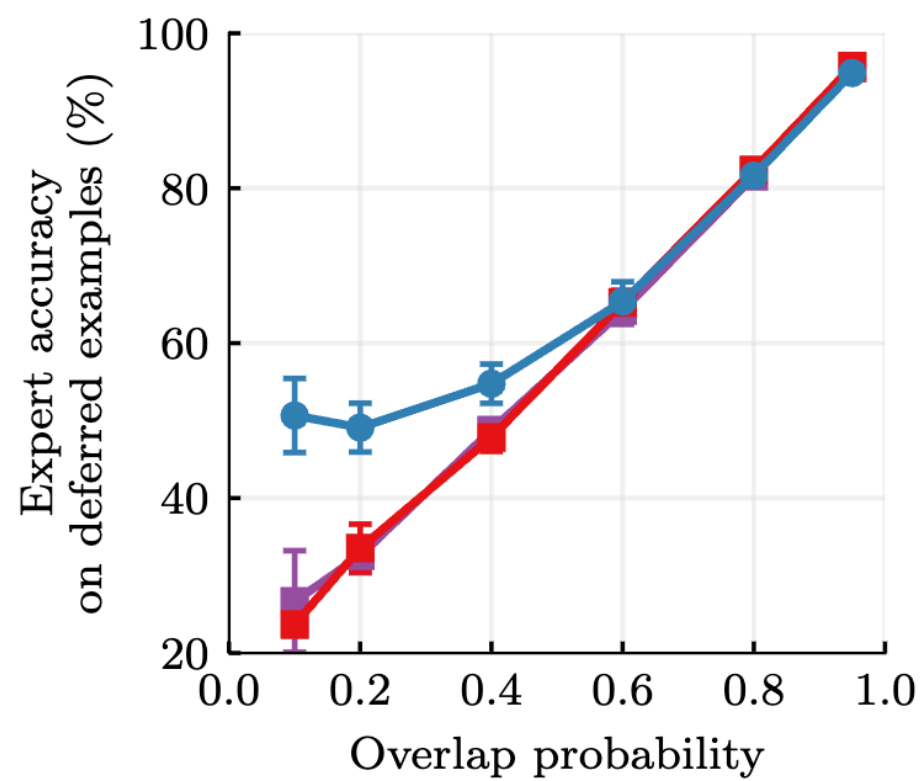
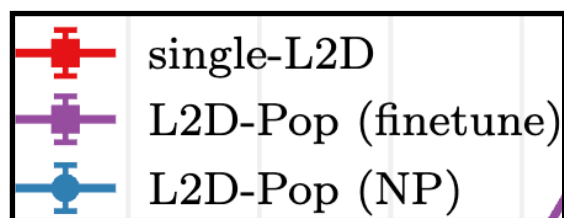
(a) Traffic Signs



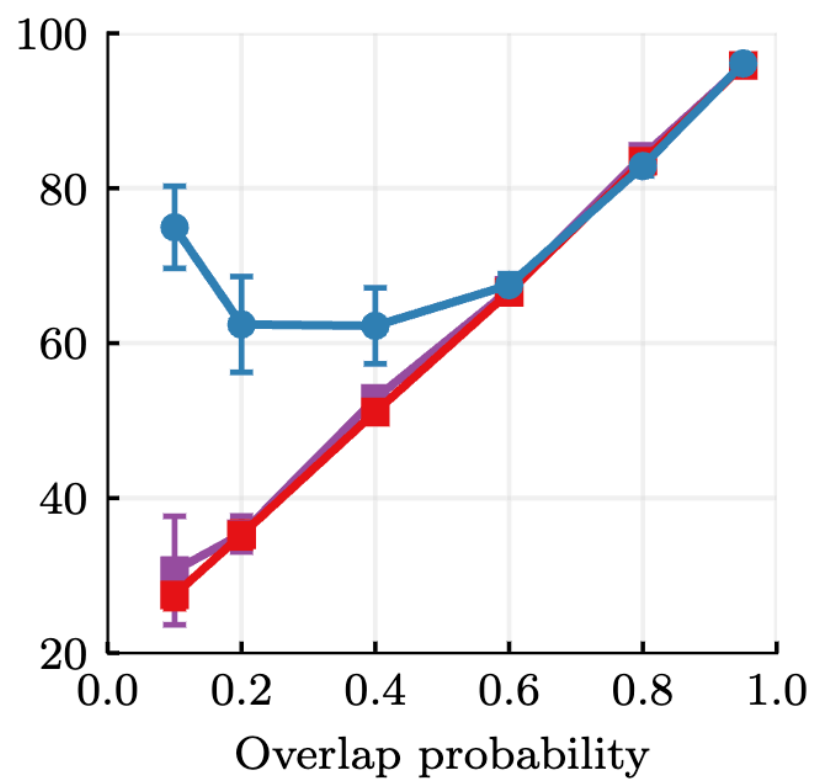
(b) CIFAR-10



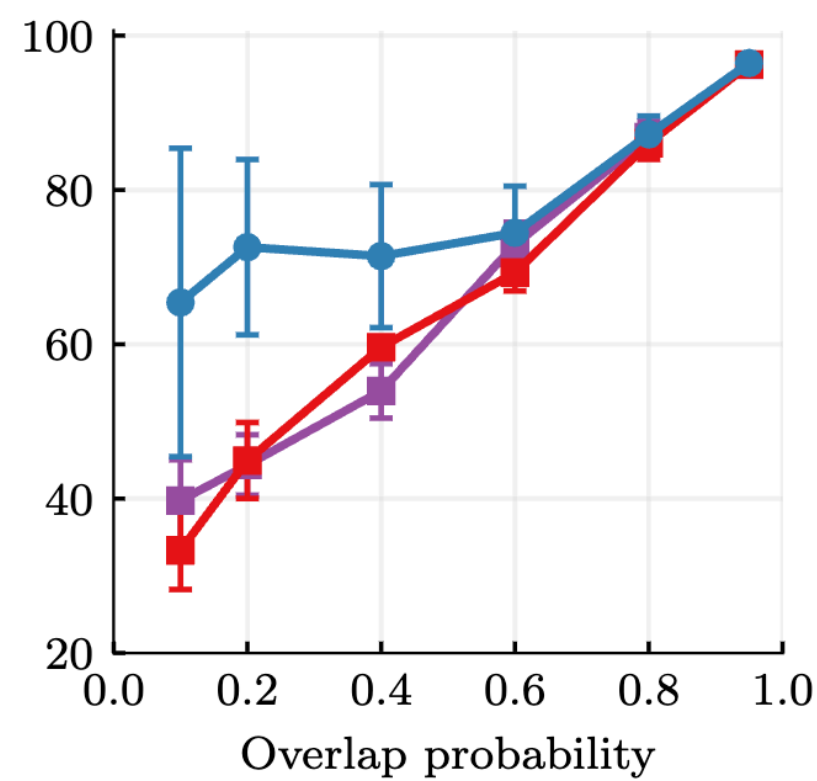
(c) HAM10000



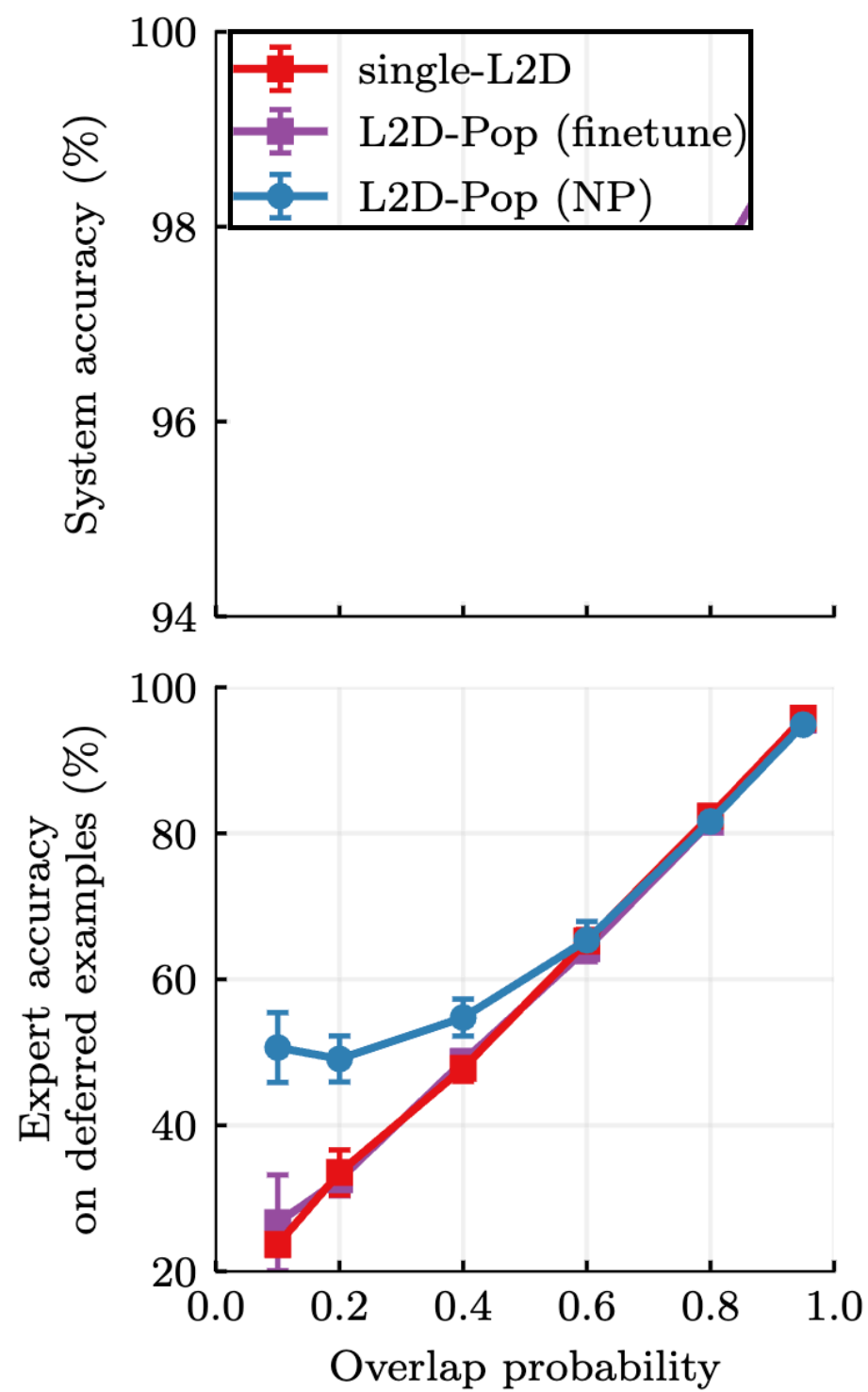
(a) Traffic Signs



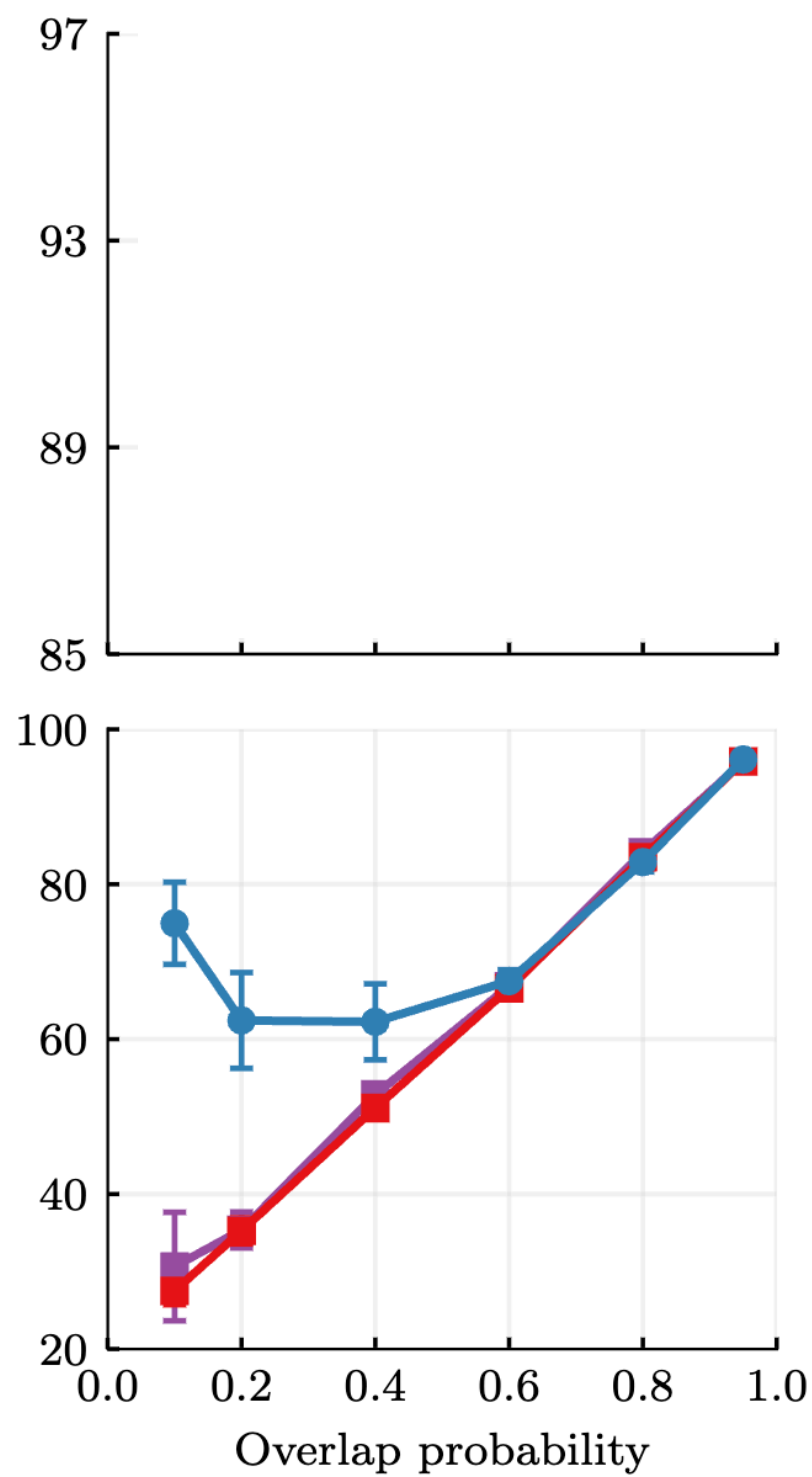
(b) CIFAR-10



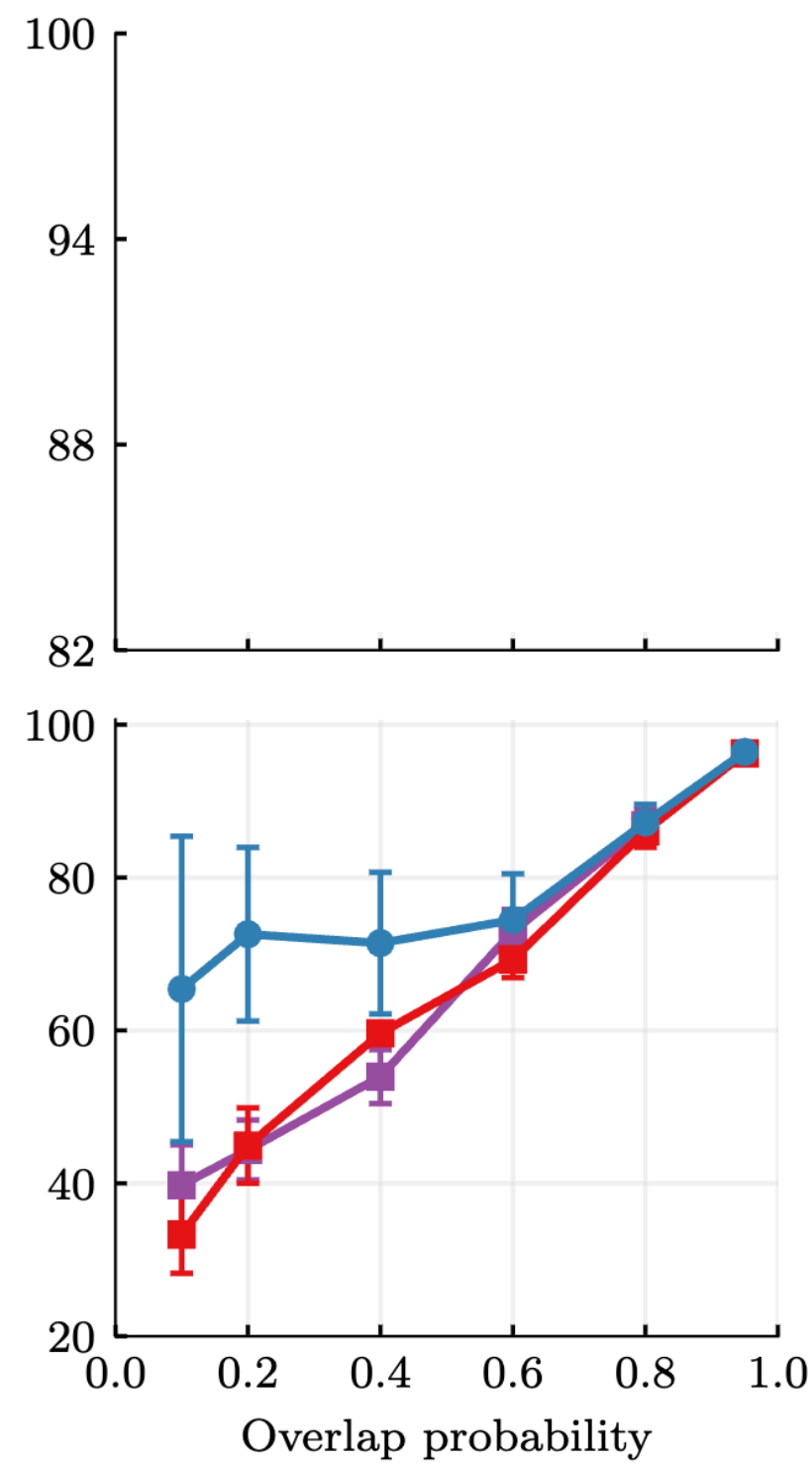
(c) HAM10000



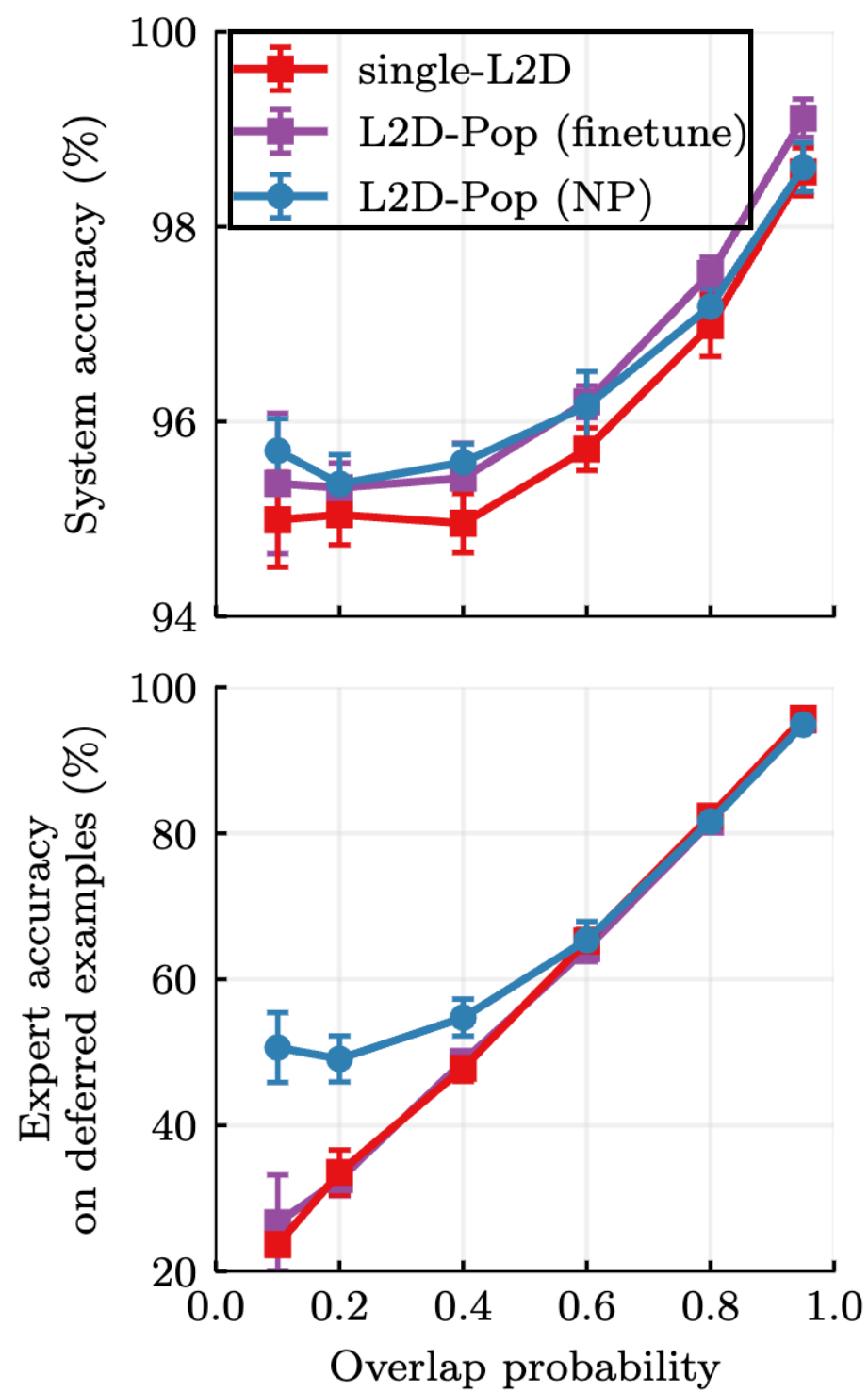
(a) Traffic Signs



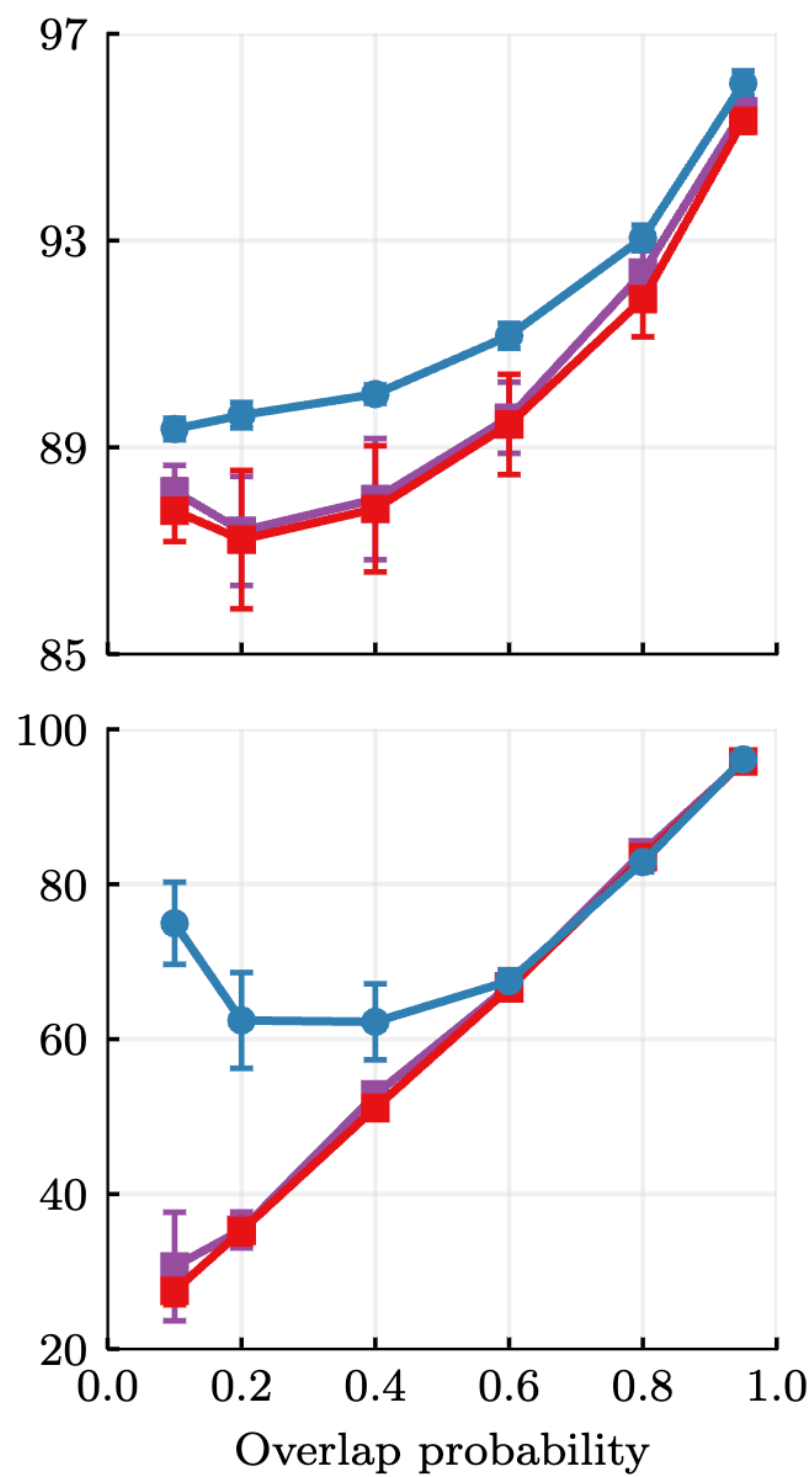
(b) CIFAR-10



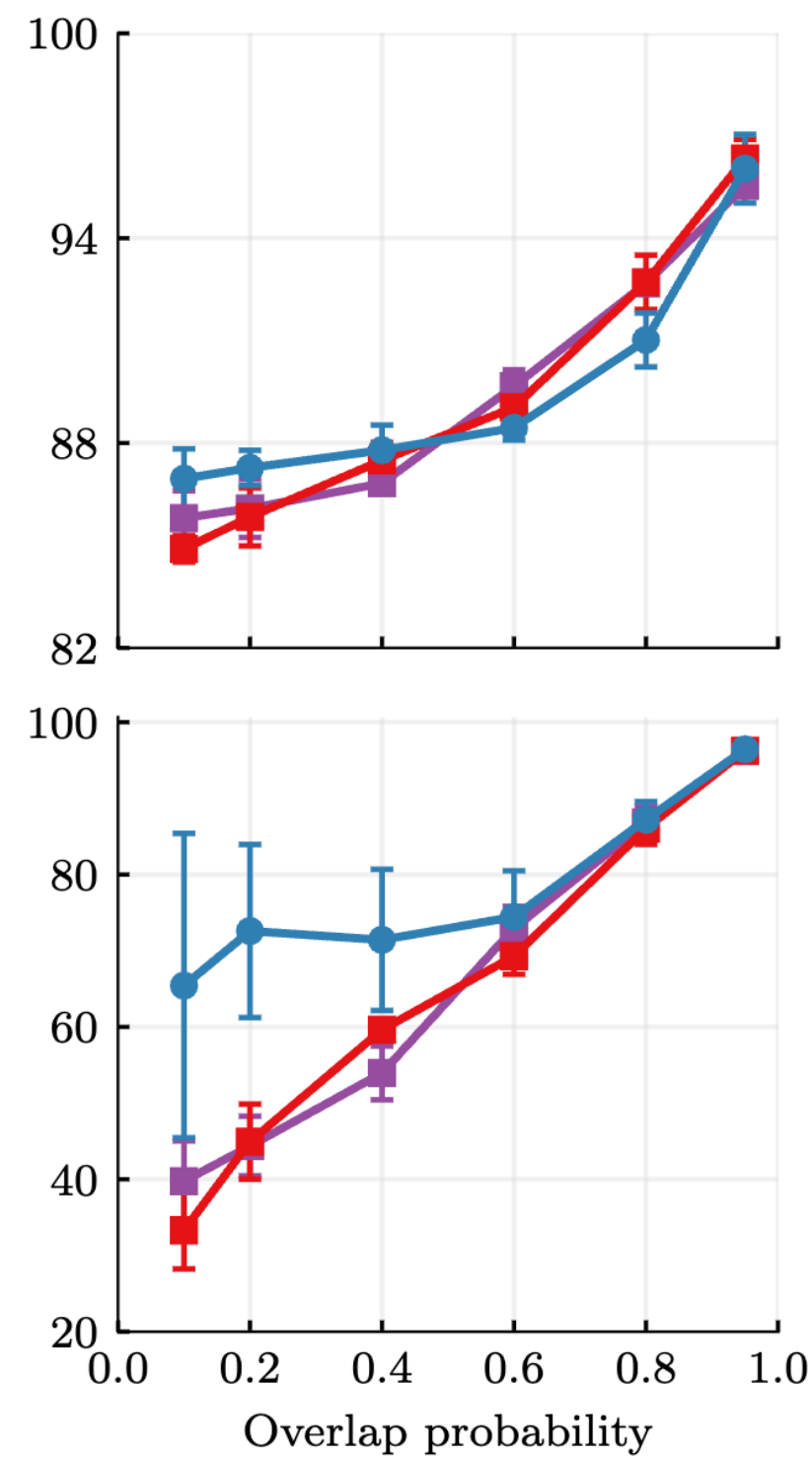
(c) HAM10000



(a) Traffic Signs



(b) CIFAR-10



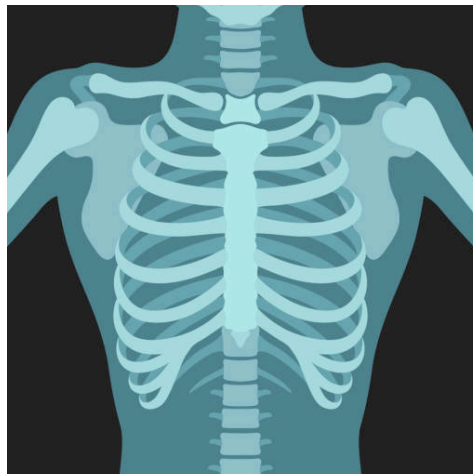
(c) HAM10000



- ⊗ single expert
  - ⊗ softmax surrogate loss
  - ⊗ improving calibration via one-vs-all
- ⊗ multiple experts
  - ⊗ surrogate losses
  - ⊗ conformal sets of experts
- ⊗ **population of experts**
  - ⊗ surrogate losses
  - ⊗ meta-learning a rejector

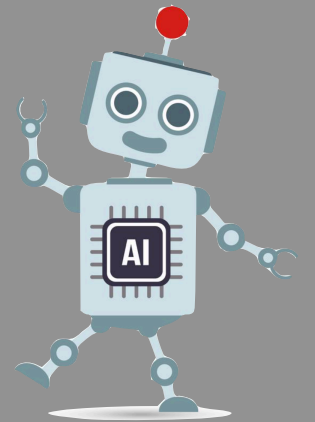
- ⊗ **single expert**
  - ⊗ softmax surrogate loss
  - ⊗ improving calibration via one-vs-all
  
- ⊗ **multiple experts**
  - ⊗ surrogate losses
  - ⊗ conformal sets of experts
  
- ⊗ **population of experts**
  - ⊗ surrogate losses
  - ⊗ meta-learning a rejector

input  
features



allocation  
mechanism

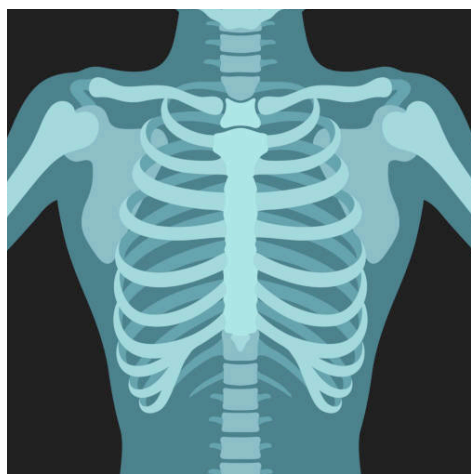
classifier



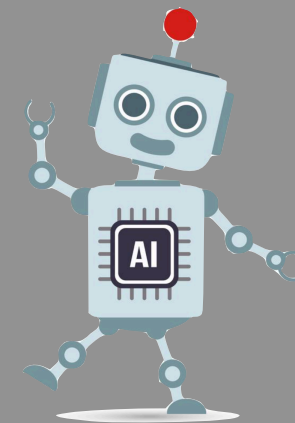
expert



input  
features



allocation  
mechanism



classifier



expert #1

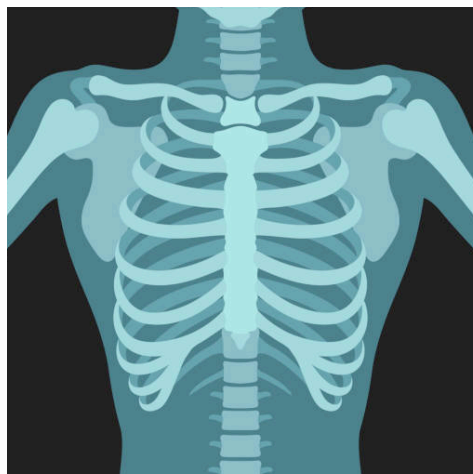


expert #3

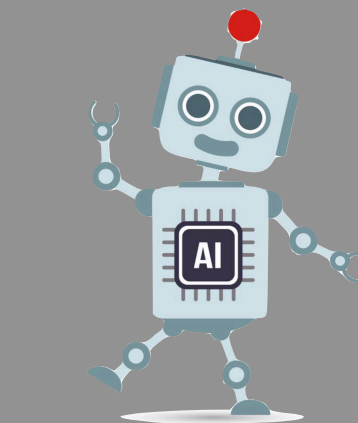


expert #2

input  
features

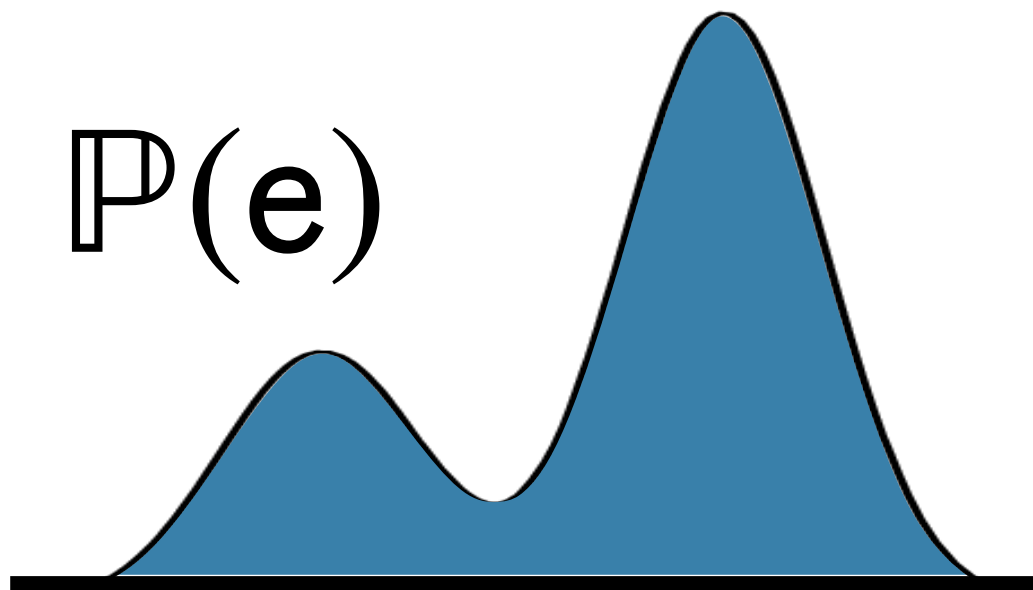


allocation  
mechanism



classifier

$\mathbb{P}(e)$



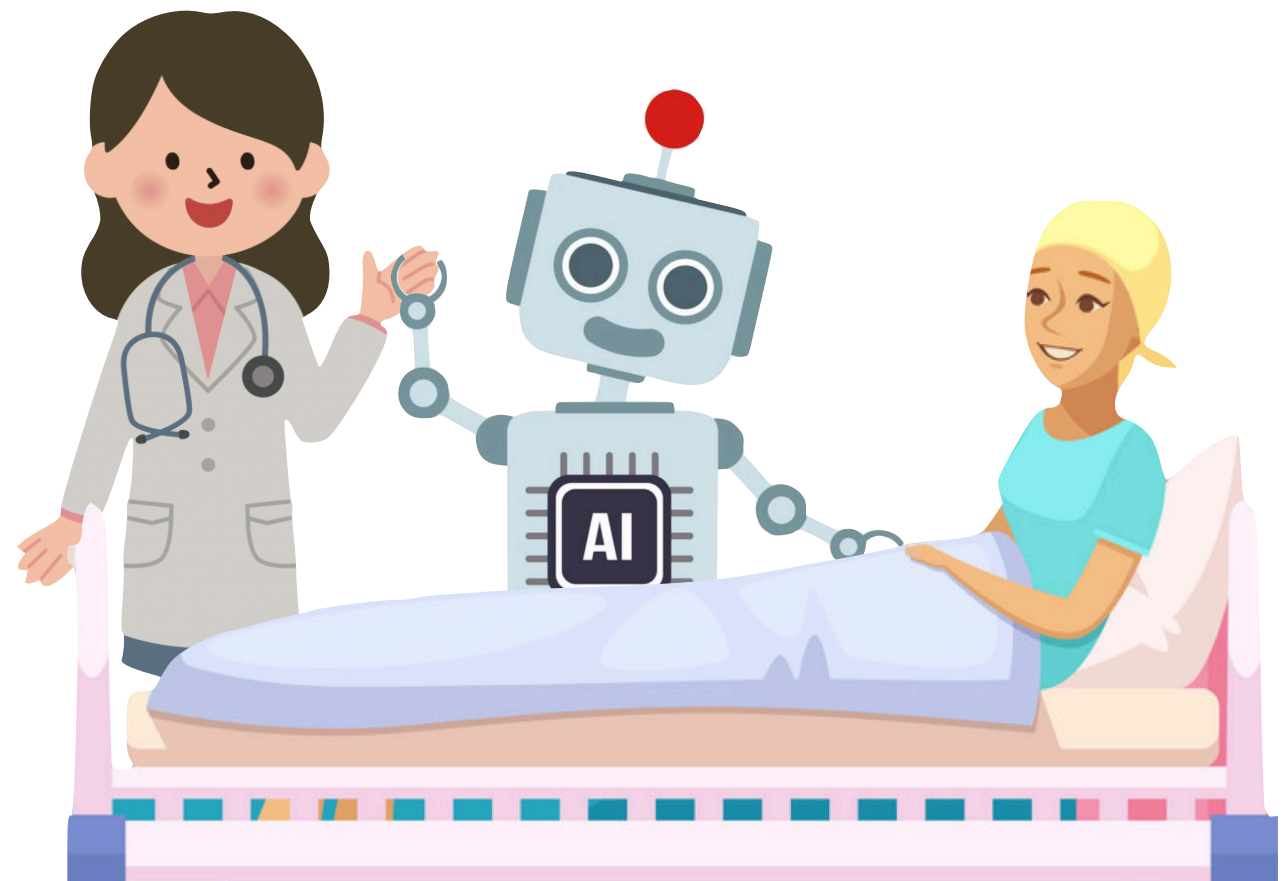
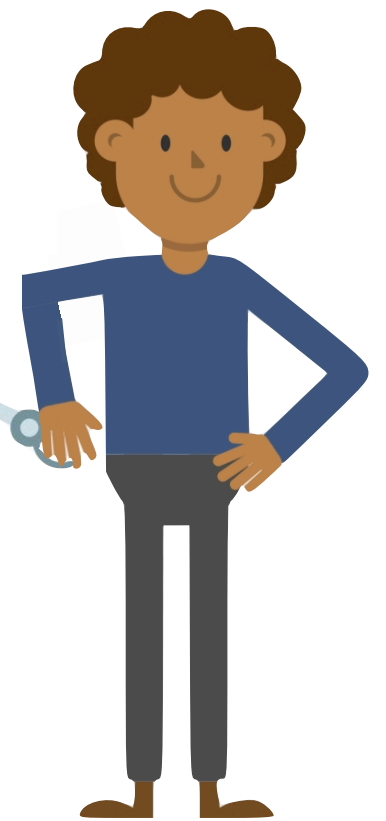
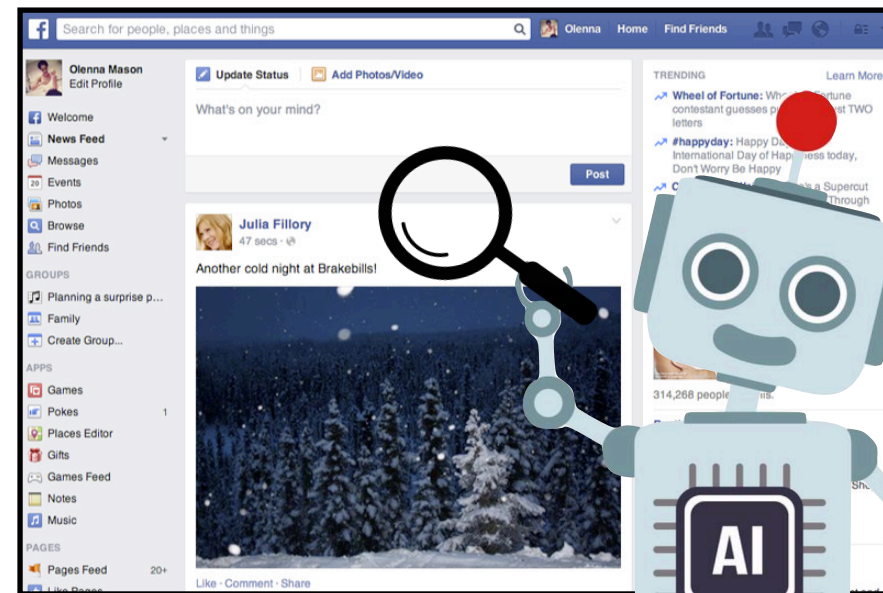
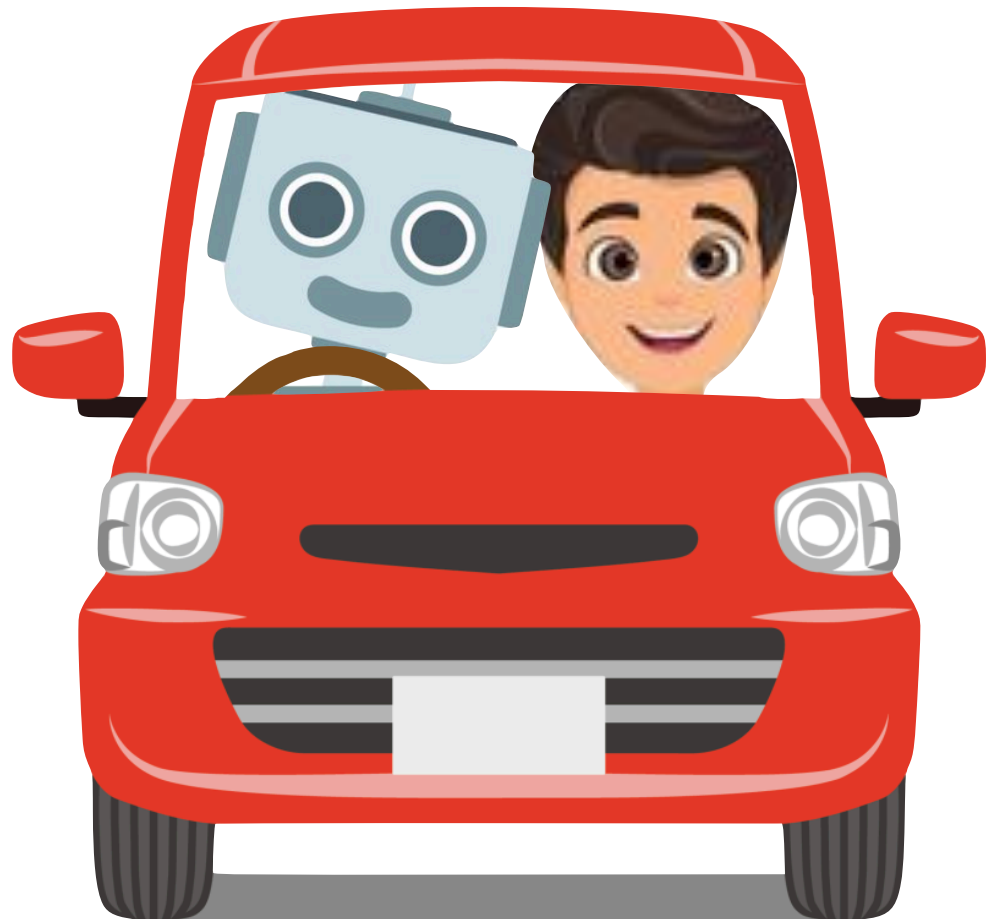
experts

$e \sim \mathbb{P}(e)$



?

expert



papers & code



funding provided by



co-authors



Rajeev  
Verma



Daniel  
Barrejón



Dharmesh  
Tailor



Putra  
Manggala



Aditya  
Patra

# Appendix



# 0-1 loss

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

$$\sum_n (1 - r(x_n)) \underbrace{\mathbb{I}[h(x_n) \neq y_n]}_{\text{classifier loss}} + r(x_n) \underbrace{\mathbb{I}[m_n \neq y_n]}_{\text{expert loss}}$$

# estimators

single expert

softmax:  $\hat{p}(m = y | \mathbf{x}) = \frac{h_{\perp}(\mathbf{x})}{1 - h_{\perp}(\mathbf{x})}$

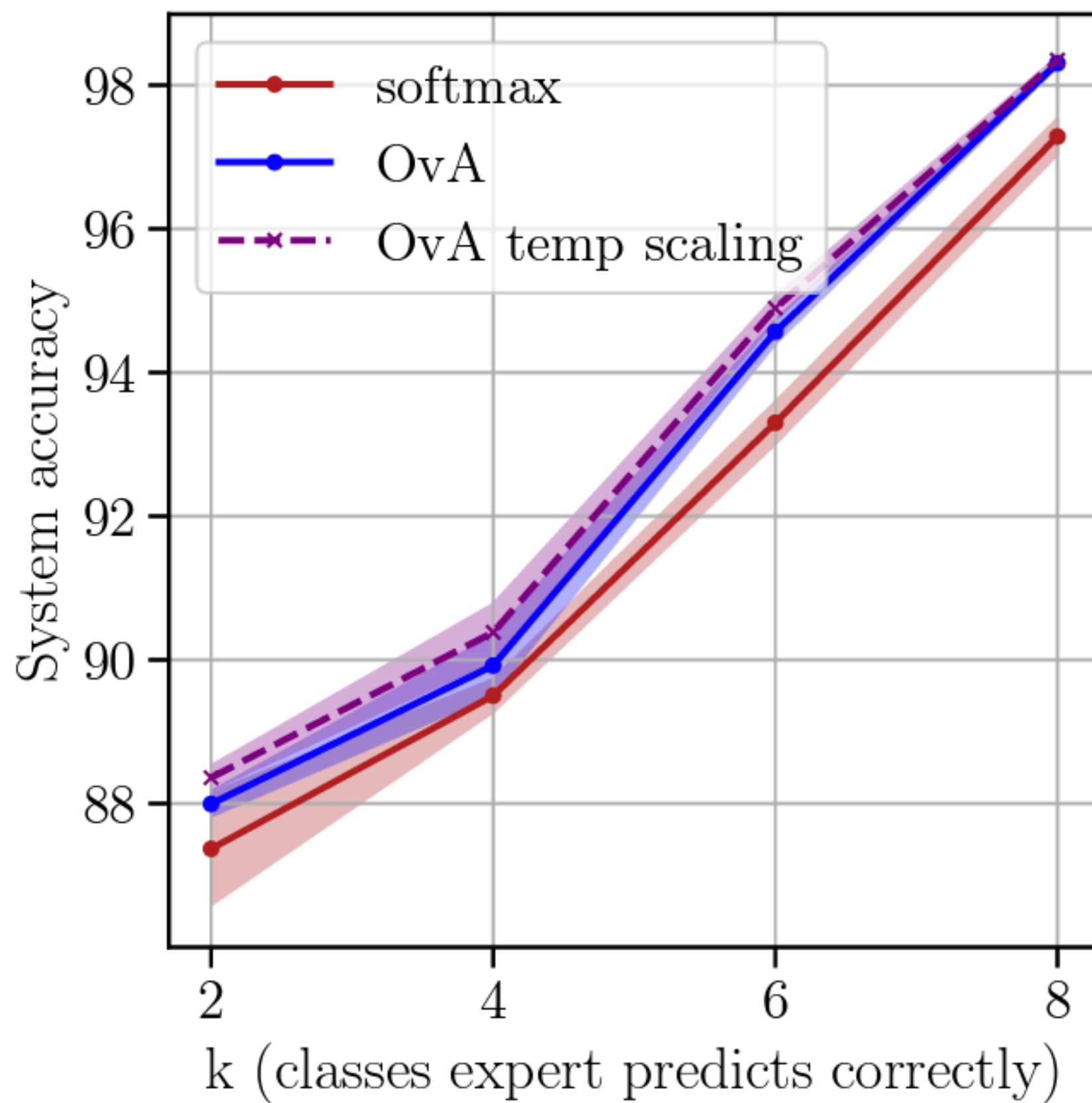
one-vs-all:  $\hat{p}(m = y | \mathbf{x}) = h_{\perp}(\mathbf{x})$

---

multi-expert

softmax:  $\hat{p}(m_j = y | \mathbf{x}) = \frac{h_{\perp,j}(\mathbf{x})}{1 - \sum_{e=1}^J h_{\perp,e}(\mathbf{x})}$

one-vs-all:  $\hat{p}(m_j = y | \mathbf{x}) = h_{\perp,j}(\mathbf{x})$

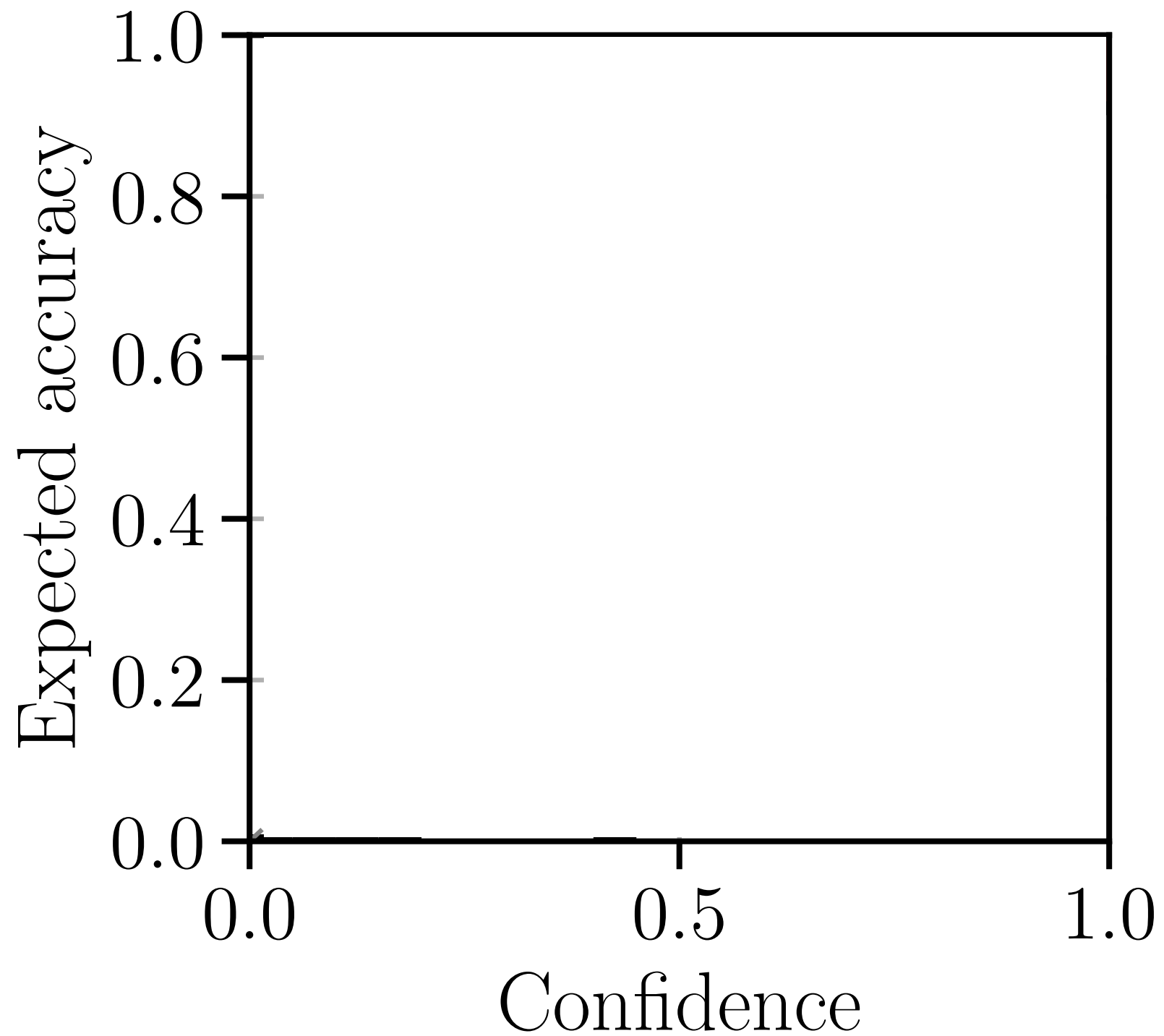


# hate speech detection

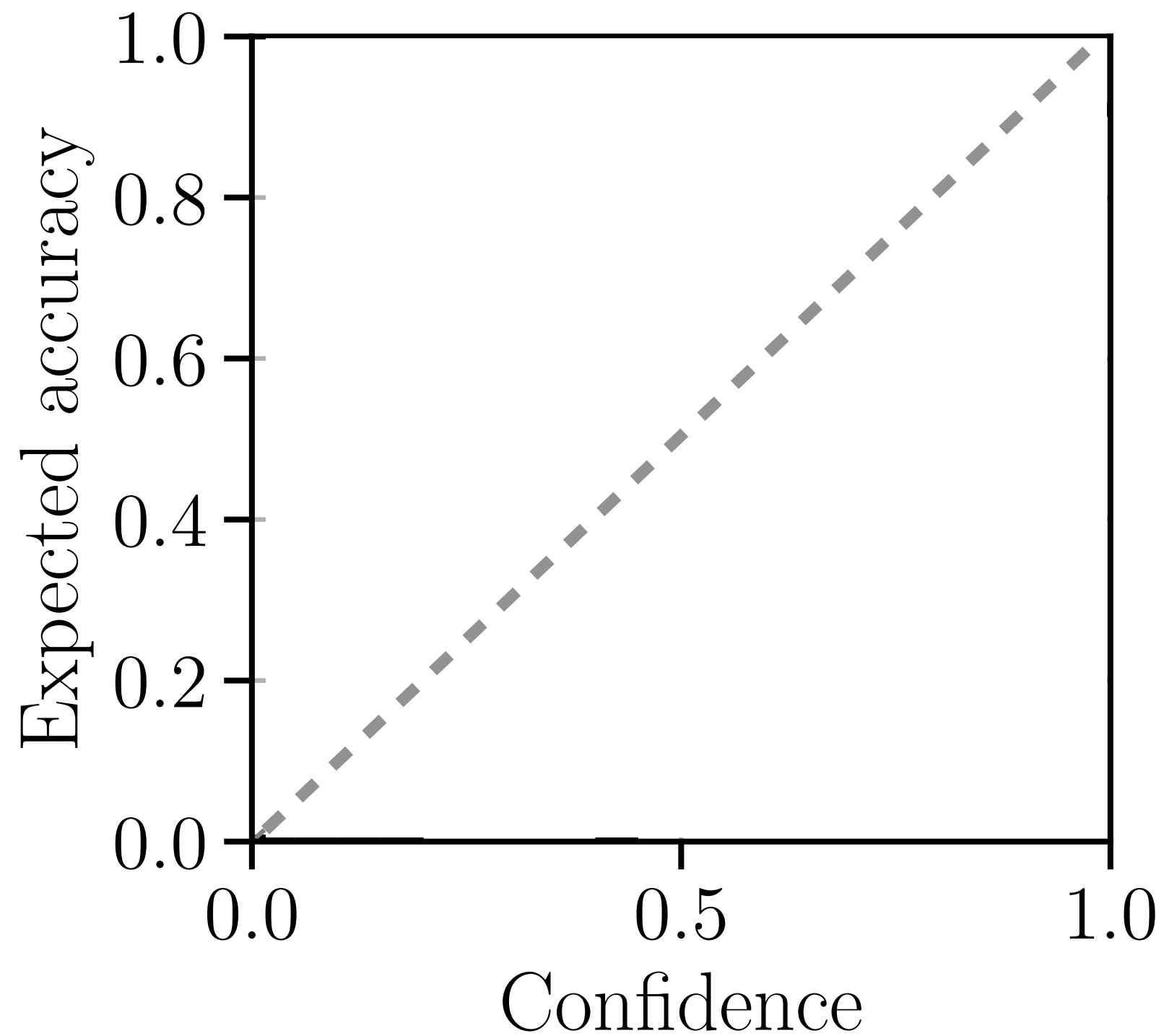


[Davidson et al., ICWSM 2017]

# hate speech detection



# hate speech detection



# hate speech detection

softmax

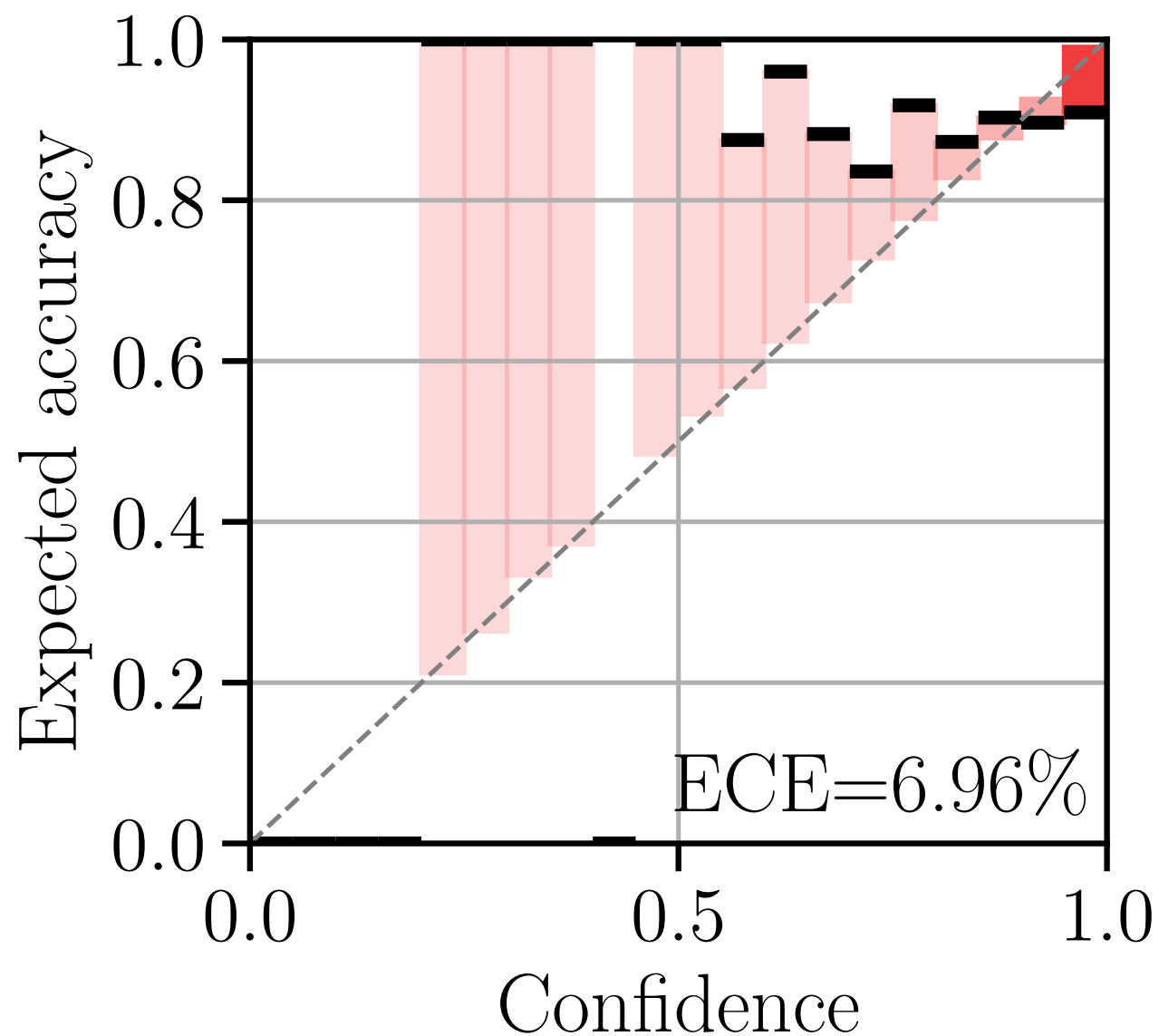
one-vs-all (ours)



# hate speech detection

softmax

one-vs-all (ours)

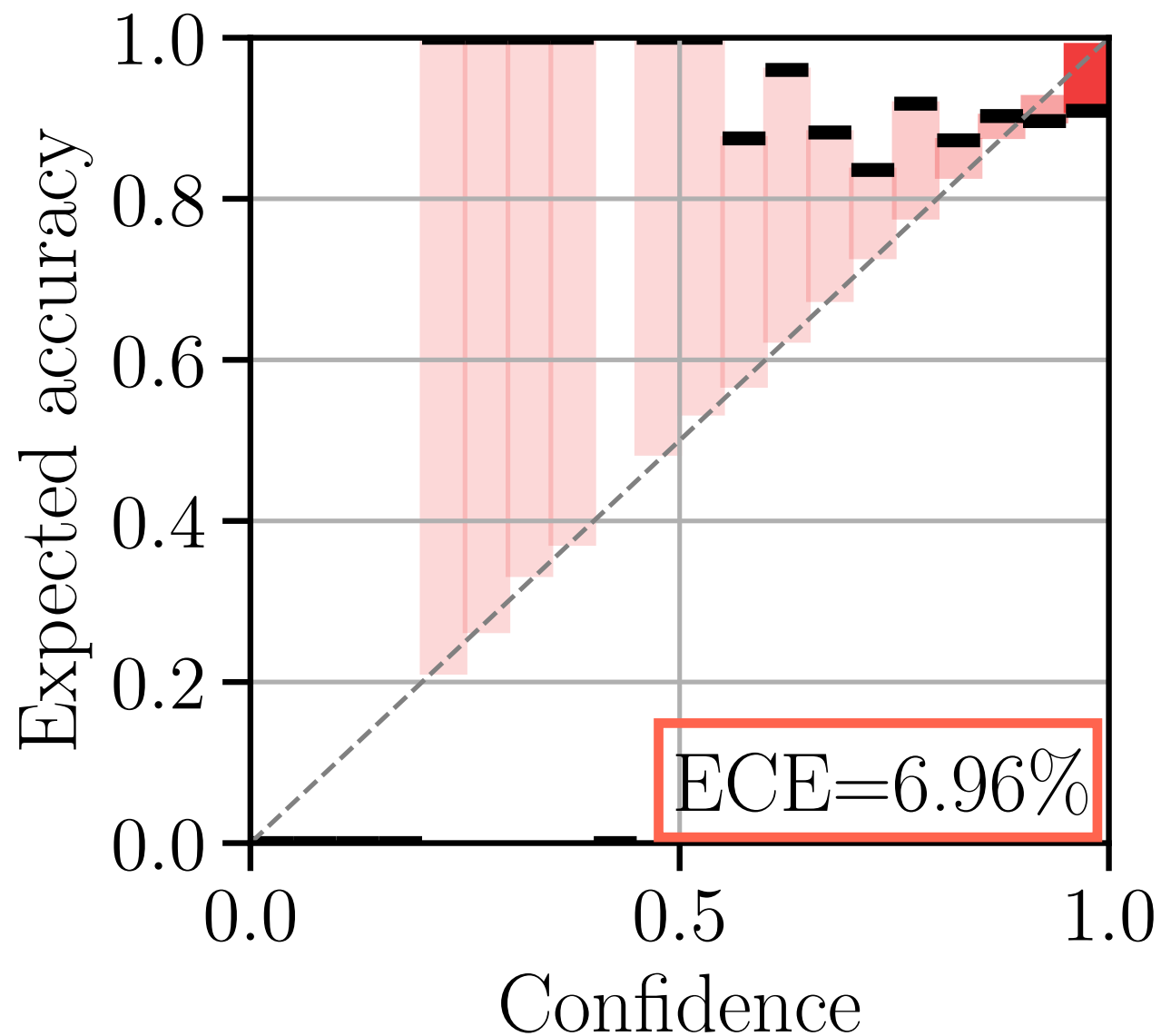




# hate speech detection

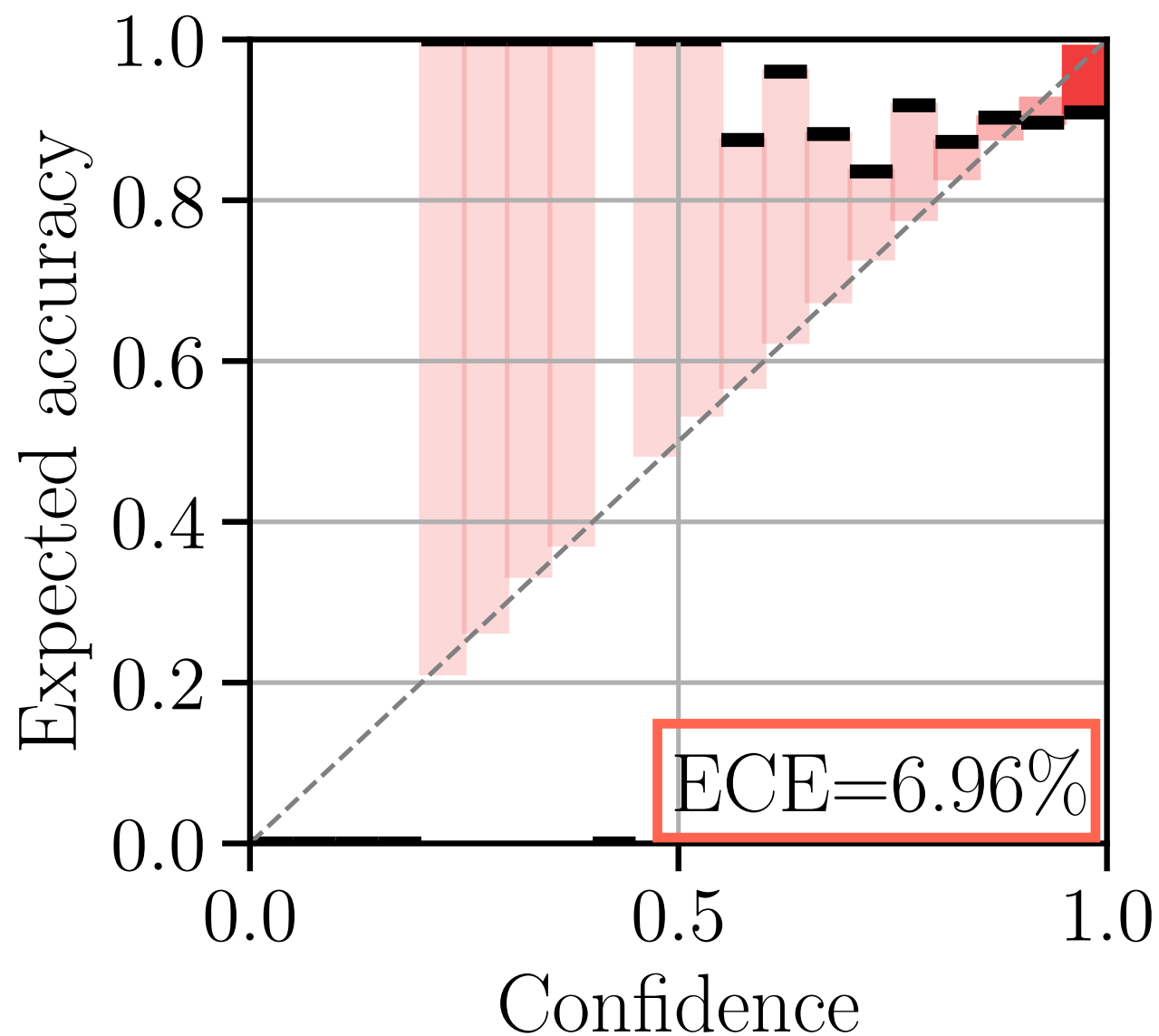
softmax

one-vs-all (ours)

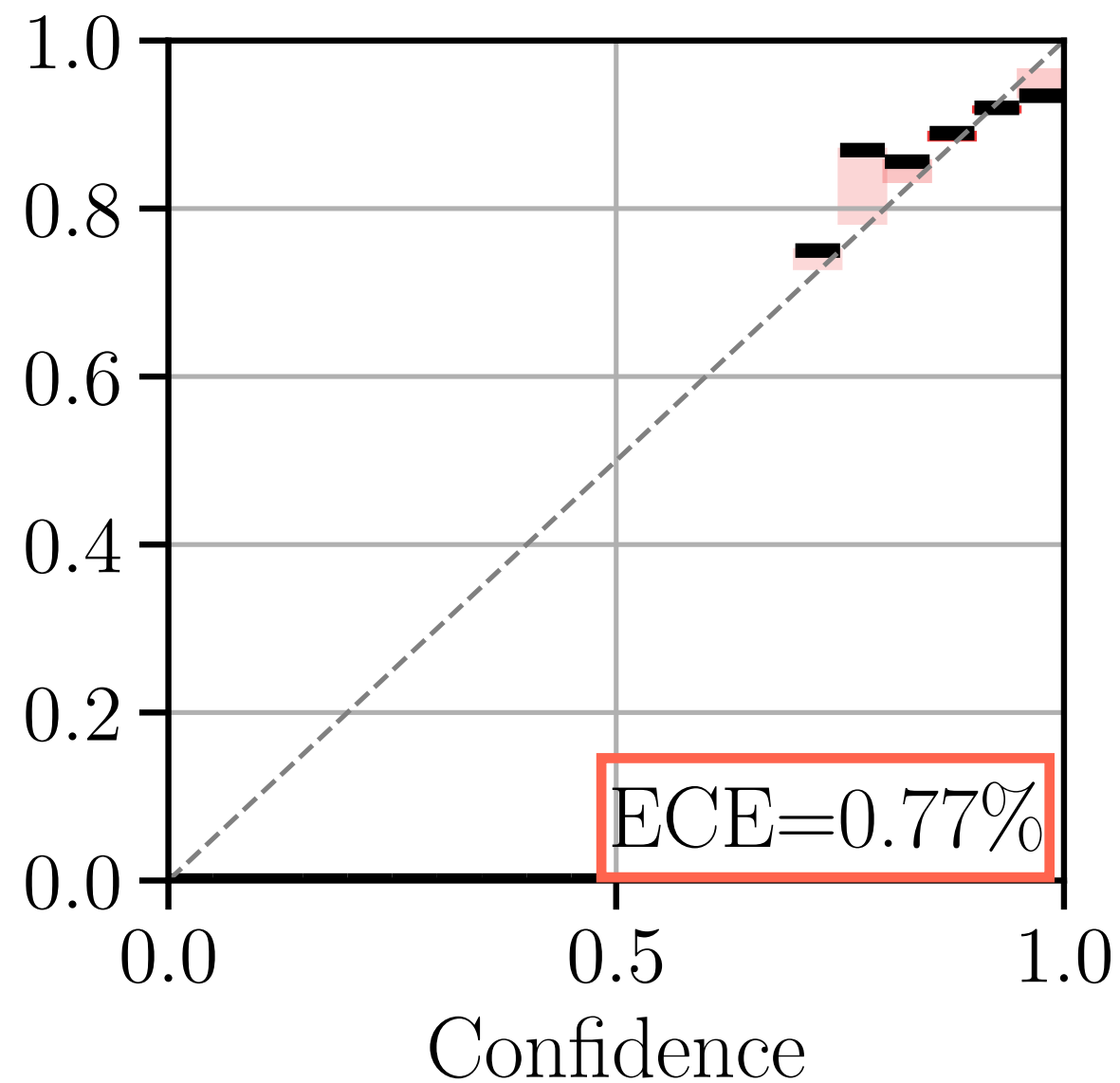


# hate speech detection

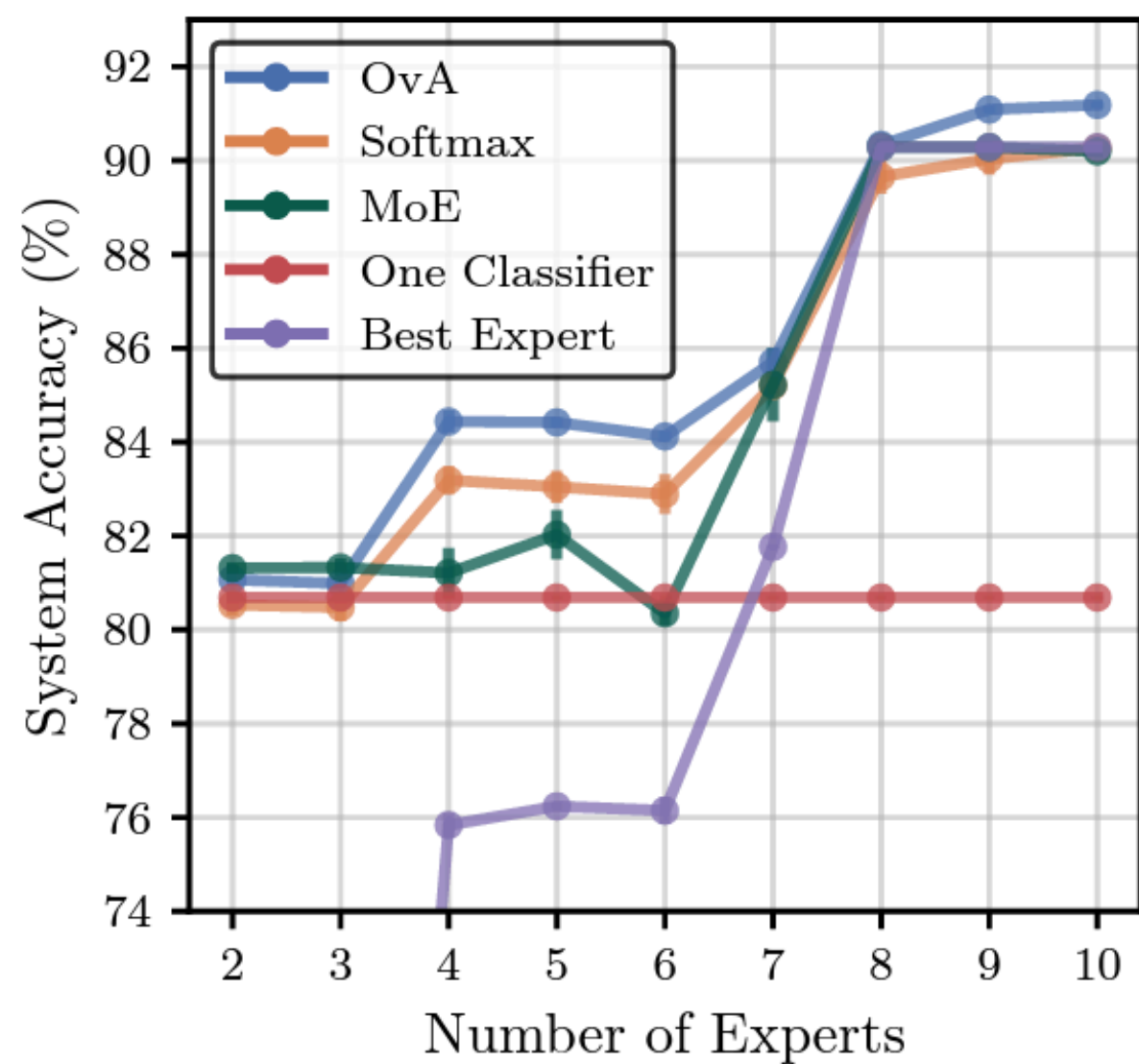
## softmax



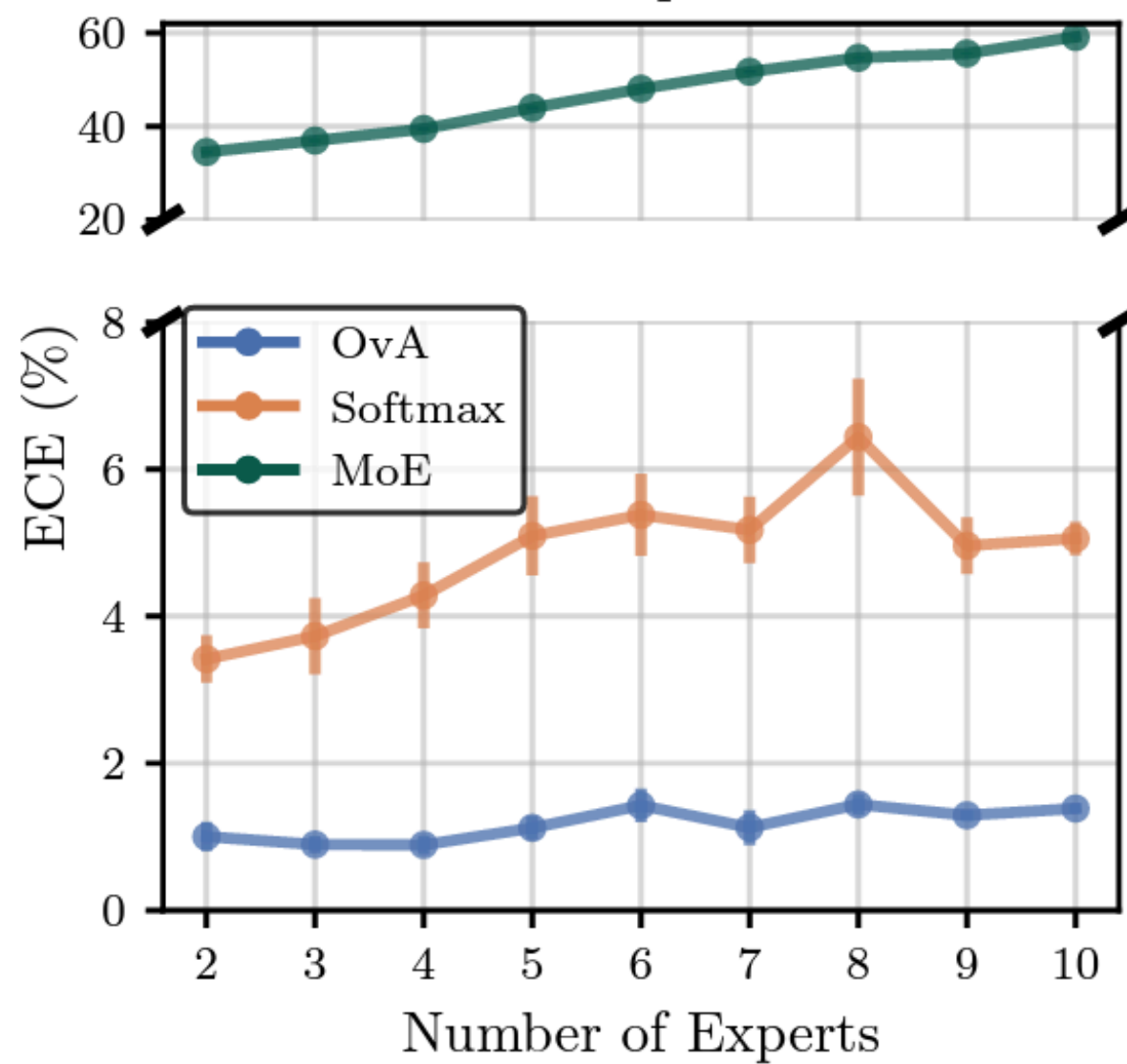
## one-vs-all (ours)



Hate Speech

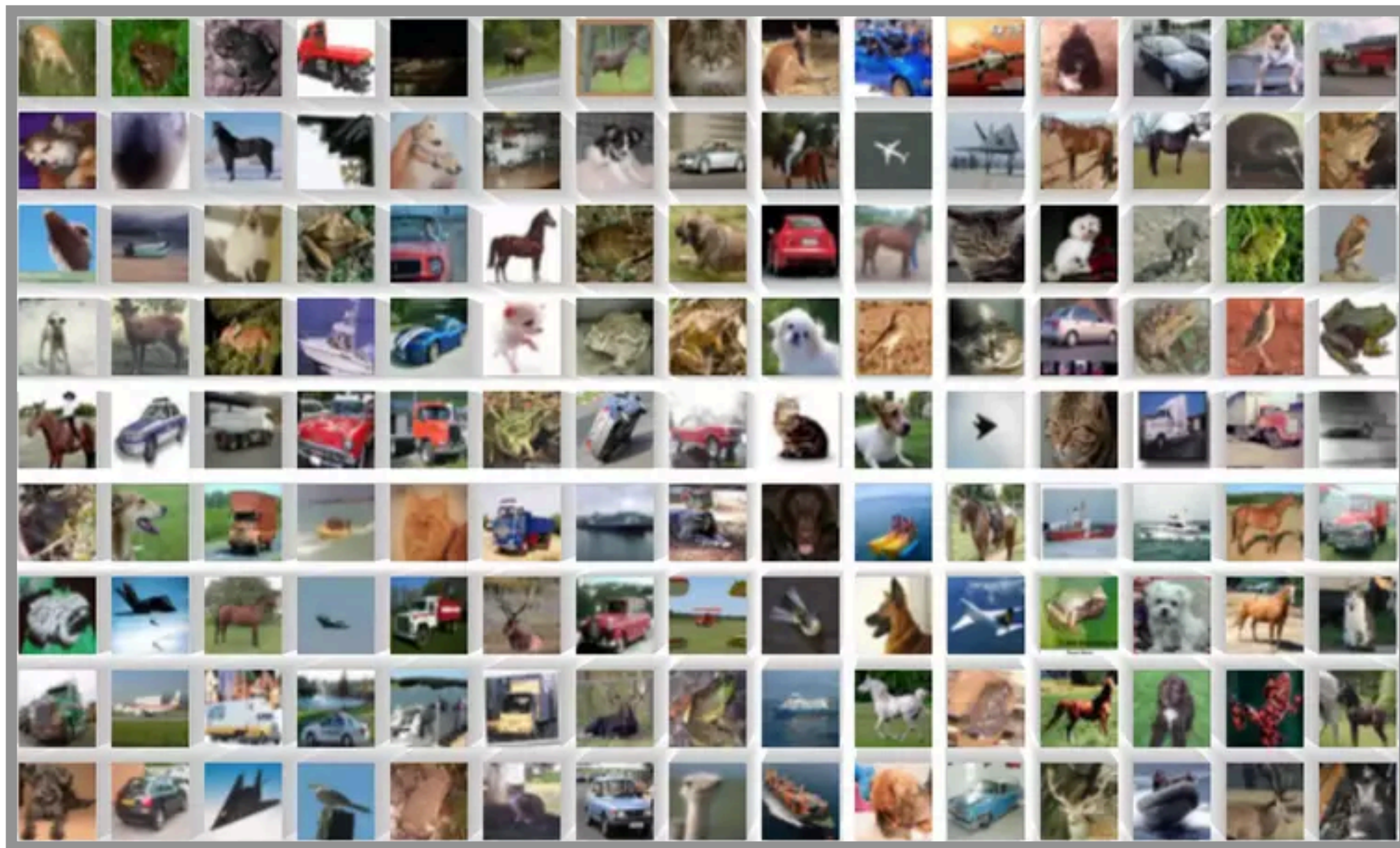


Hate Speech

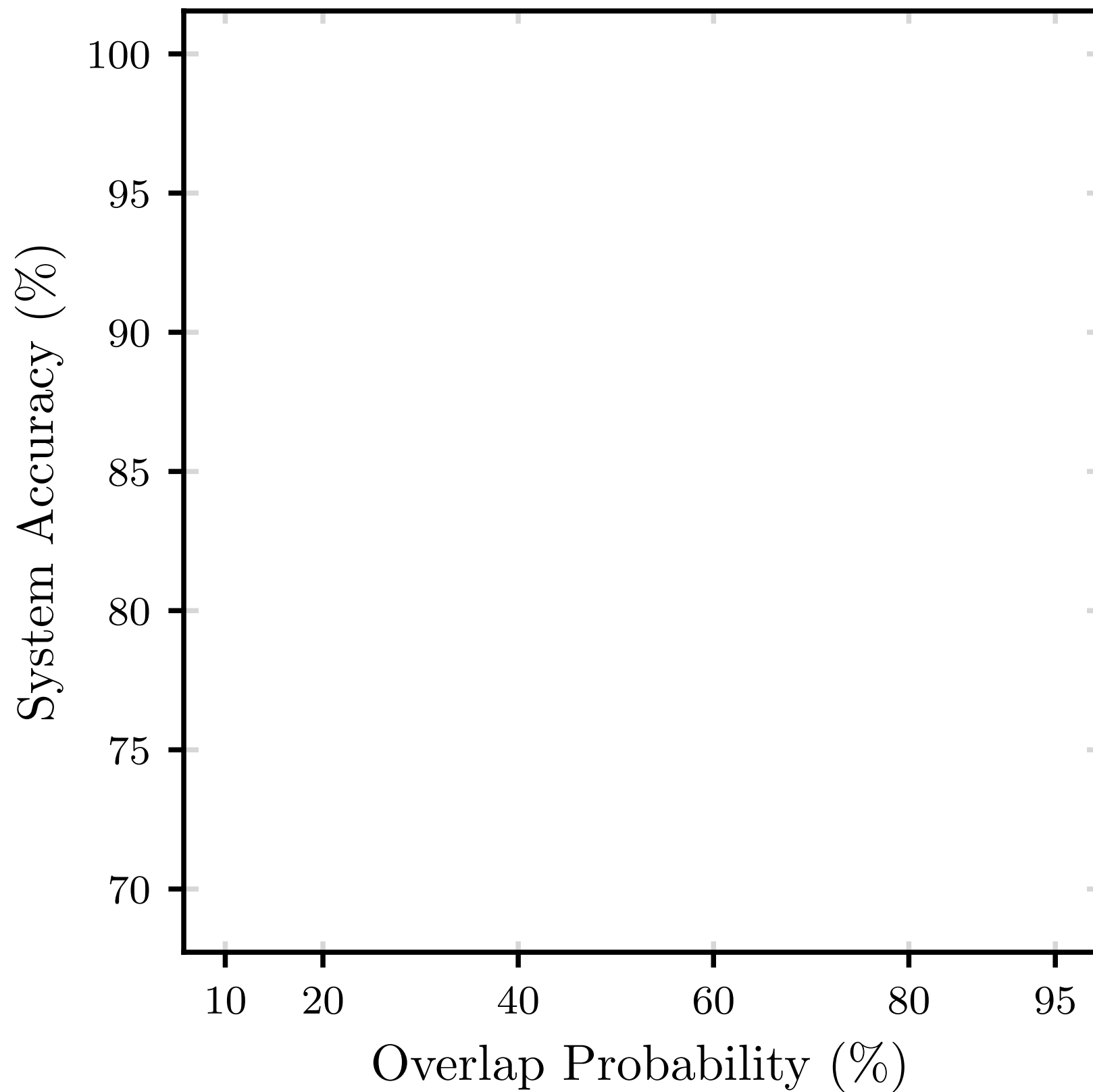


conformal: downstream performance

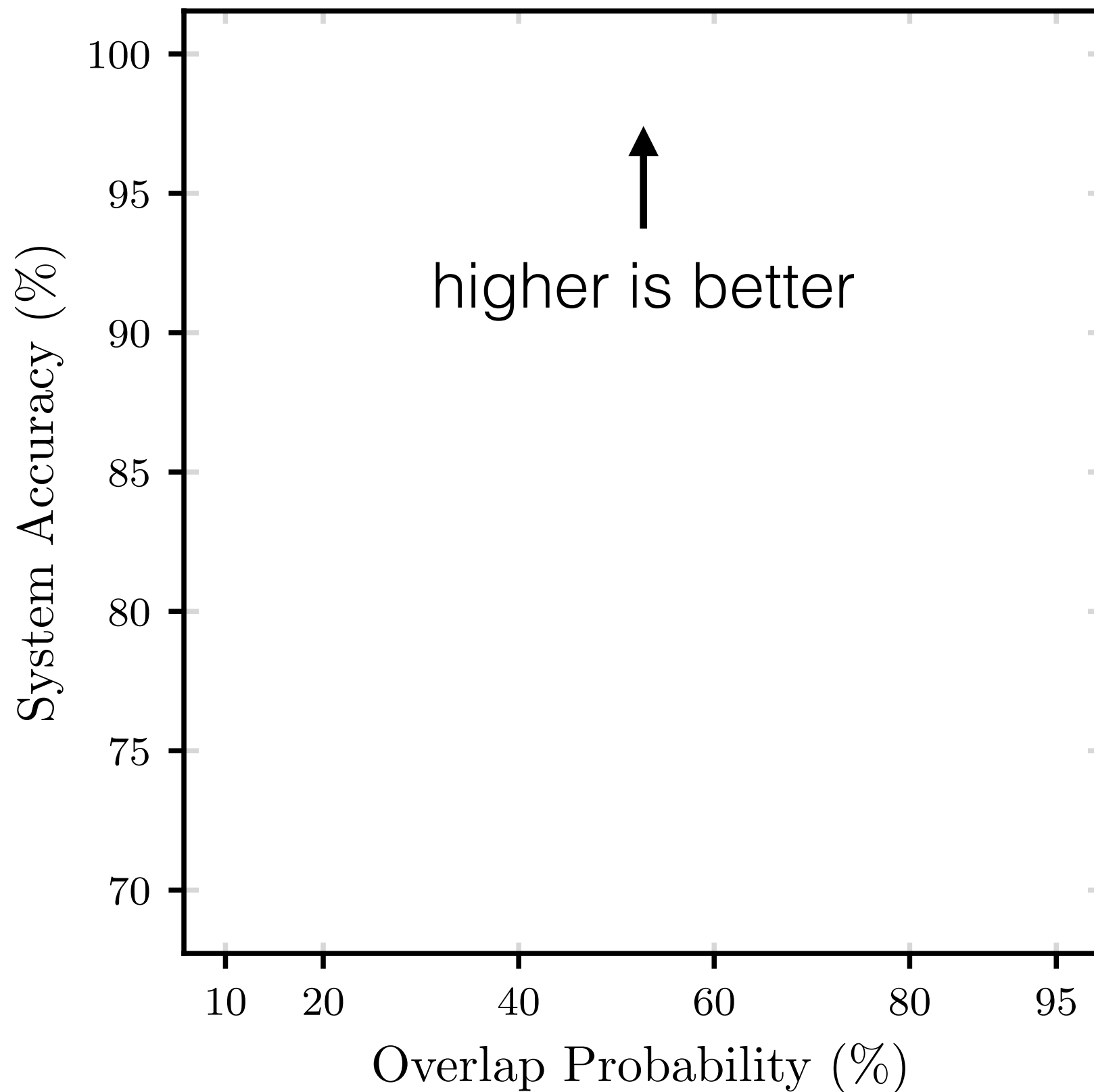
CIFAR-10



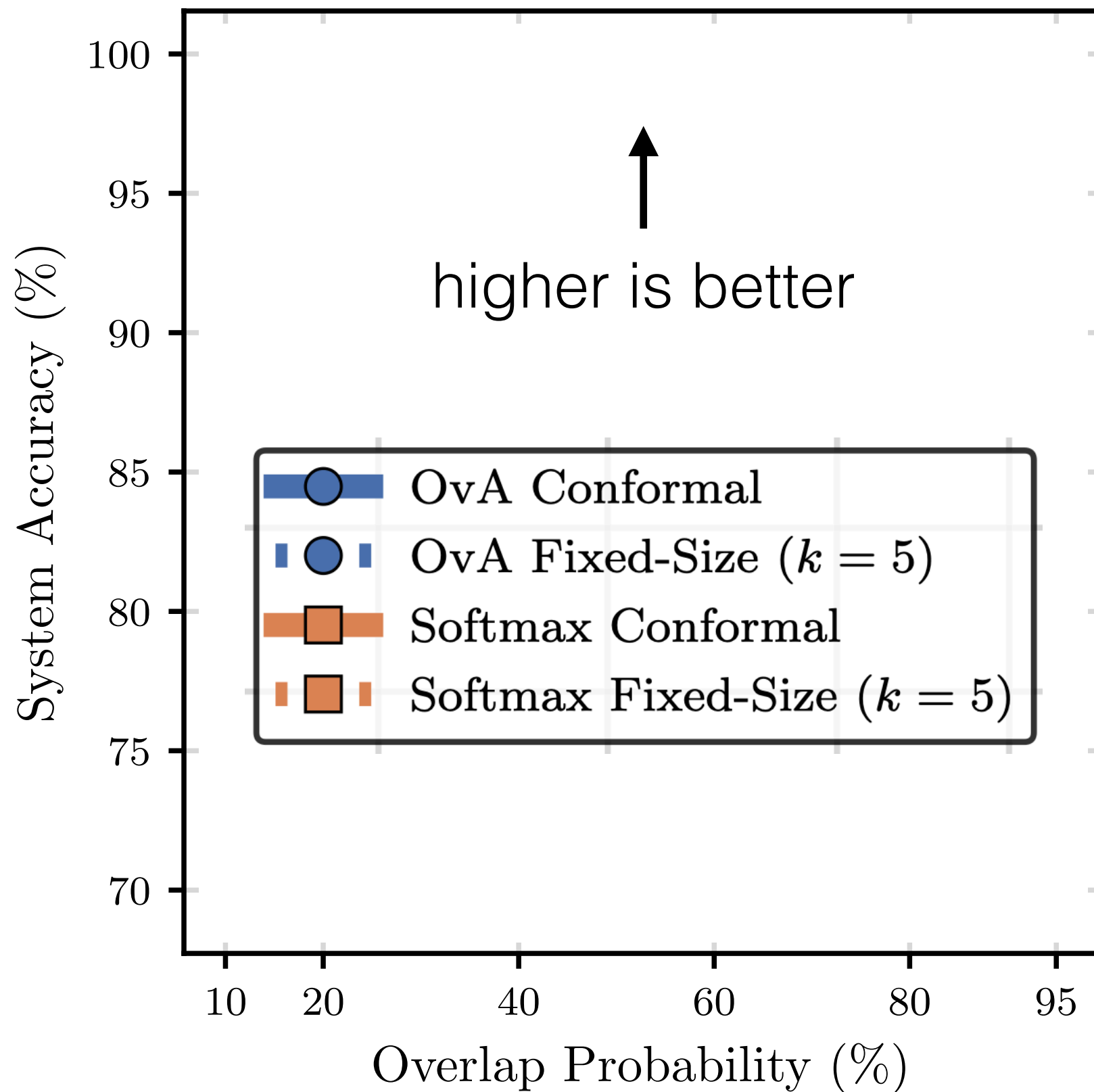
# downstream performance



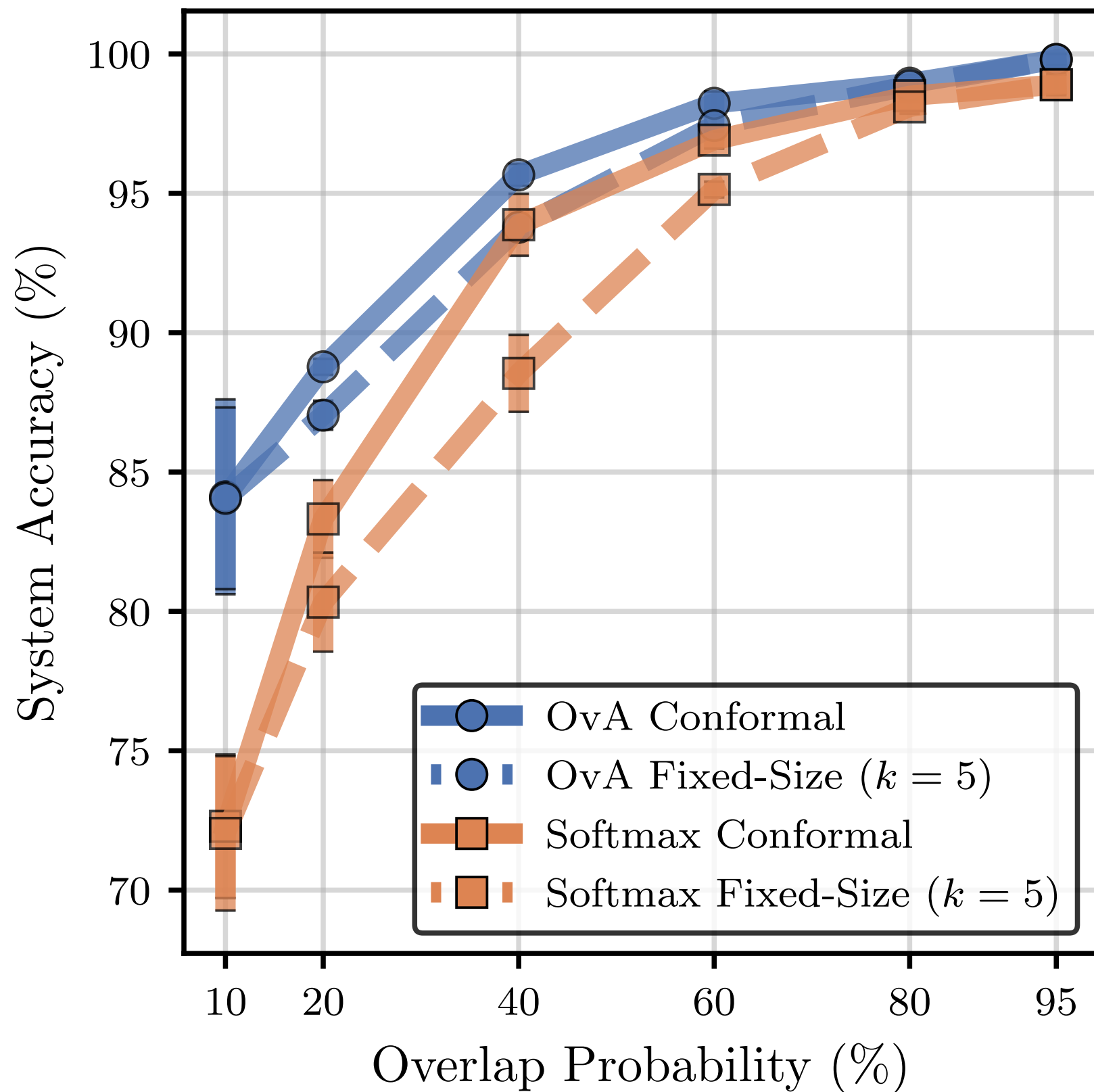
# downstream performance



# downstream performance



# downstream performance





# simulated experts:

Table 2: HAM10000 experts configuration.

	Expert configuration	$p_{in}$ [%]	$p_{out}$ [%]	Diagnostic Category [in]
1	Random Expert	-	-	[nv, bkl, df, vasc, mel, bcc, akiec]
2	Dermatologist for malign	25	15	[mel, bcc, akiec]
3	Dermatologist for benign	25	15	[nv, bkl, df, vasc]
4	Specialized dermatologist in nv	50	15	[nv]
5	Specialized dermatologist in vasc	70	15	[vasc]
6	Specialized dermatologist in mel	75	15	[mel]
7	Dermatologist for benign	75	25	[nv, bkl, df, vasc]
8	MLP Mixer	-	-	[nv, bkl, df, vasc, mel, bcc, akiec]
9	Experienced dermatologist	80	50	[nv, bkl, df, vasc, mel, bcc, akiec]
10	Experienced dermatologist	80	60	[nv, bkl, df, vasc, mel, bcc, akiec]

# simulated experts:

Table 1: Hate Speech and Galaxy-Zoo experts configuration.

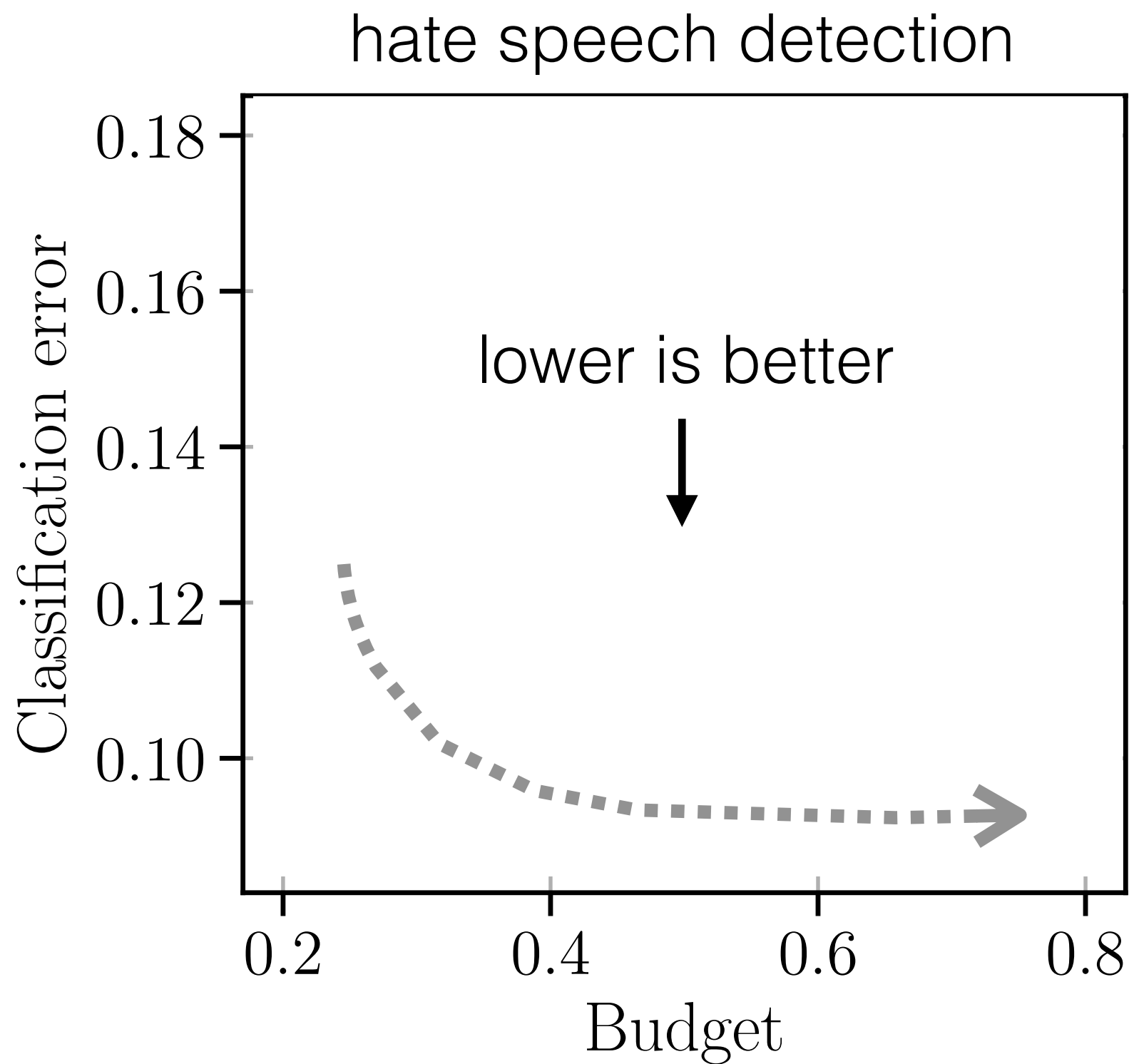
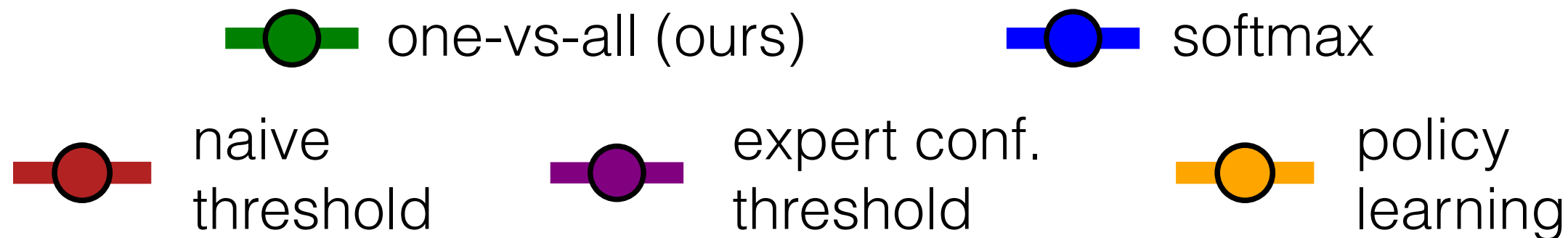
	Expert configuration	$p_{\text{flip}}[\%]$	$p_{\text{annotator}}[\%]$
1	Random Expert	-	-
2	Probabilistic Expert	-	10
3	Flipping Human Expert	50	-
4	Probabilistic Expert	-	75
5	Flipping Human Expert	30	-
6	Flipping Human Expert	20	-
7	Probabilistic Expert	-	85
8	Human Expert	-	-
9	Probabilistic Expert	-	50
10	Human Expert	-	-



# hate speech detection

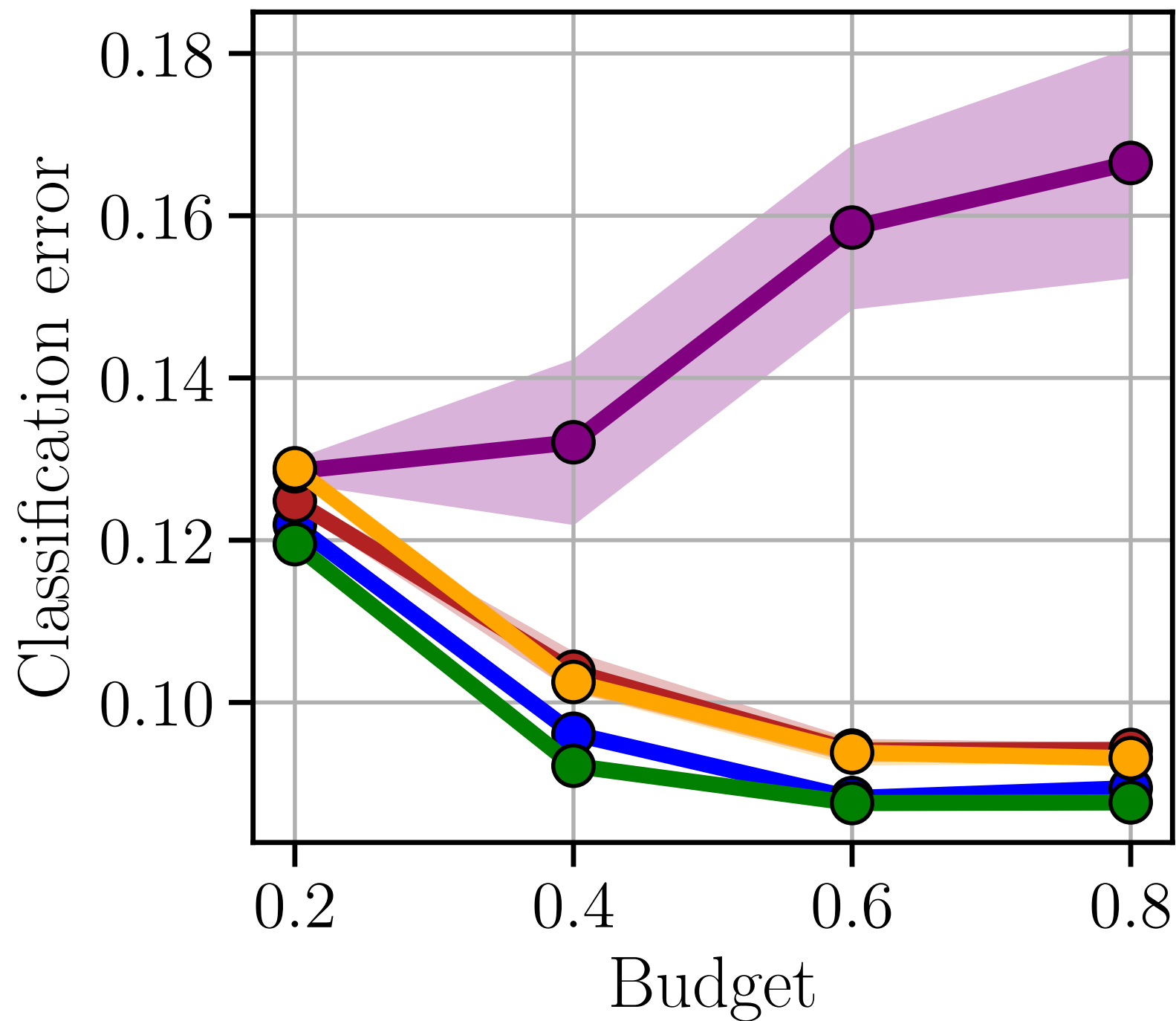


[Davidson et al., ICWSM 2017]





### hate speech detection



# conformal inference: train-time



expert #1



expert #2



expert #3

# conformal inference: train-time



expert #1

$$h_{\perp,1}(x)$$



expert #2

$$h_{\perp,2}(x)$$



expert #3

$$h_{\perp,3}(x)$$

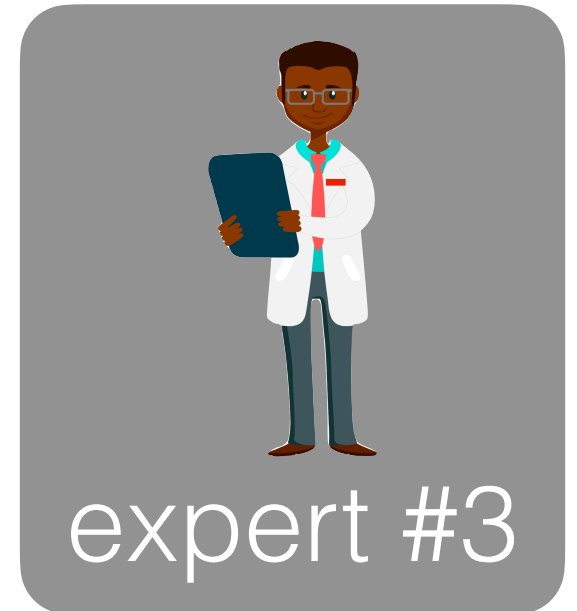
# conformal inference: train-time



$$h_{\perp,1}(x)$$



$$h_{\perp,2}(x)$$



$$h_{\perp,3}(x)$$

using validation data, compute the  
(1- $\alpha$ )-quantile of a conformity statistic:

$$\hat{q}_{1-\alpha}$$



# conformal inference: test-time



expert #1

$$h_{\perp,1}(x)$$



expert #2

$$h_{\perp,2}(x)$$



expert #3

$$h_{\perp,3}(x)$$

# conformal inference: test-time



$$h_{\perp,3}(x) > h_{\perp,1}(x) > h_{\perp,2}(x)$$

# conformal inference: test-time



$$h_{\perp,3}(x) > h_{\perp,1}(x) > h_{\perp,2}(x)$$

$$C(x) = \left\{ \begin{array}{l} \text{check if:} \\ \sum_{e \in C(x)} h_{\perp,e}(x) \geq \hat{q}_{1-\alpha} \end{array} \right. ?$$

# conformal inference: test-time



$h_{\perp,3}(x)$

$>$



expert #1

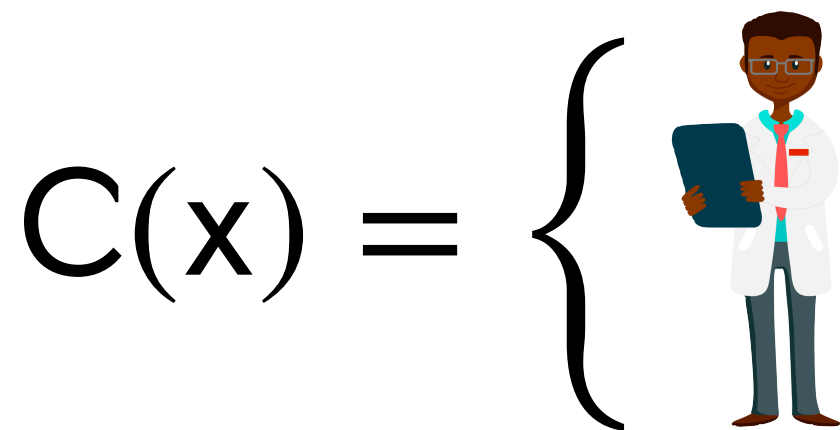
$h_{\perp,1}(x)$

$>$



expert #2

$h_{\perp,2}(x)$



$= \left\{ \right.$

check if:

$h_{\perp,3} \stackrel{?}{\geq} \hat{q}_{1-\alpha}$

# conformal inference: test-time



$h_{\perp,3}(x)$

$>$



expert #1

$h_{\perp,1}(x)$

$>$



expert #2

$h_{\perp,2}(x)$

$$C(x) = \left\{ \text{expert #3} \right\}$$

check if:

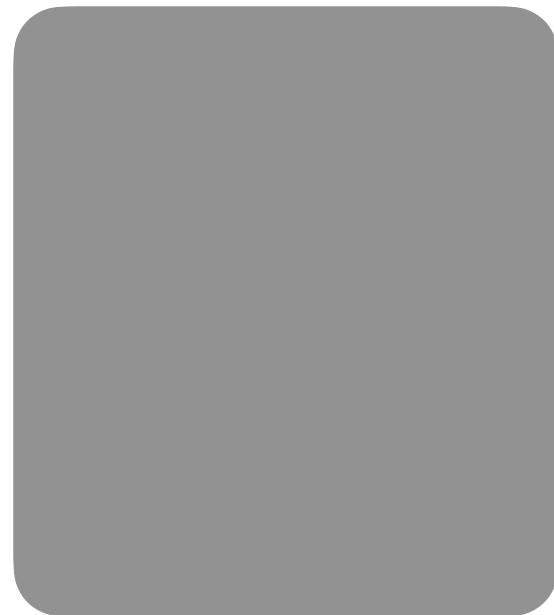
$h_{\perp,3} \overset{\text{X}}{\geq} \hat{q}_{1-\alpha}$

# conformal inference: test-time



$h_{\perp,3}(x)$

$>$



$h_{\perp,1}(x)$

$>$



$h_{\perp,2}(x)$

$C(x) = \left\{ \begin{array}{c} \text{male doctor} \\ \text{female doctor} \end{array} \right\}$

check if:

$$h_{\perp,3} + h_{\perp,1} \stackrel{?}{\geq} \hat{q}_{1-\alpha}$$

# conformal inference: test-time



$h_{\perp,3}(x)$

$>$



$h_{\perp,1}(x)$

$>$



$h_{\perp,2}(x)$

$$C(x) = \left\{ \text{male doctor}, \text{female doctor} \right\}$$

check if:

$$h_{\perp,3} + h_{\perp,1} \stackrel{\checkmark}{\geq} \hat{q}_{1-\alpha}$$

Estimating  $\mathbb{P}(m = y | x)$

